

Applied Discrete Structures

Algebraic Structures Chapters 11-16

Alan Doerr and Kenneth Levasseur
Department Of Mathematical Sciences
University of Massachusetts Lowell

Version 1.0
March 2012



Home

Blog

Errata

Home: <http://faculty.uml.edu/klevasseur/ADS2/>
Blog: <http://applieddiscretestructures.blogspot.com/>
Errata: <http://faculty.uml.edu/klevasseur/ADS2/errata.html>



Applied Discrete Structures by Alan Doerr & Kenneth Levasseur is licensed under a Creative Commons Attribution-Noncommercial-ShareAlike 3.0 United States License.

(<http://creativecommons.org/licenses/by-nc-sa/3.0/us/>)

Previously published by Pearson Education, Inc. under the title *Applied Discrete Structures for Computer Science*

Applied Discrete Structures

To our families
Donna, Christopher, Melissa, and Patrick Doerr
and
Karen, Joseph, Kathryn, and Matthew Levasseur



Table of Contents

Preface

Chapter 11 Algebraic Systems

- 11.1 Operations**
- 11.2 Algebraic Systems**
- 11.3 Some General Properties of Groups**
- 11.4 \mathbb{Z}_n , the Integers Modulo n**
- 11.5 Subsystems**
- 11.6 Direct Products**
- 11.7 Isomorphisms**
- 11.8 Using Computers to Study Groups**
- Supplementary Exercises for Chapter 11**

Chapter 12 More Matrix Algebra

- 12.1 Systems of Linear Equations**
- 12.2 Matrix Inversion**
- 12.3 An Introduction to Vector Spaces**
- 12.4 The Diagonalization Process**
- 12.5 Some Applications**
- Supplementary Exercises for Chapter 12**

Chapter 13 Boolean Algebra

- 13.1 Posets Revisited**
- 13.2 Lattices**
- 13.3 Boolean Algebras**
- 13.4 Atoms of a Boolean Algebra**
- 13.5 Finite Boolean Algebras as n -tuples of Zeros and Ones**
- 13.6 Boolean Expressions**
- 13.7 A Brief Introduction to the Application of Boolean Algebra to Switching Theory**
- Supplementary Exercises for Chapter 13**

Chapter 14 Monoids and Automata

- 14.1 Monoids**
- 14.2 Free Monoids and Languages**
- 14.3 Automata, Finite-state Machines**
- 14.4 The Monoid of A Finite-state Machine**
- 14.5 The Machine of A Monoid**
- Supplementary Exercises for Chapter 14**

Chapter 15 Groups Theory and Applications

- 15.1 Cyclic Groups**
- 15.2 Cosets and Factor Groups**

15.3 Permutation Groups

15.4 Normal Subgroups and Group Homomorphisms

15.5 Coding Theory—Group Codes

Supplementary Exercises for Chapter 15

Chapter 16 An Introduction to Rings and Fields

16.1 Rings—Basic Definitions and Concepts

16.2 Fields

16.3 Polynomial Rings

16.4 Field Extensions

16.5 Power Series

Supplementary Exercises for Chapter 16

Solutions and Hints to Selected Exercises

Preface - what a difference 21 years make!

This is *Applied Discrete Structures, Part II - Algebraic Structures*, which contains an introduction to groups, monoids, rings, fields, vector spaces, lattices, and boolean algebras. It corresponds with the content of Discrete Structures II at UMass Lowell, which is a required course for students in Computer Science. It presumes background contained in *Part I - Fundamentals*, which is the content of Discrete Structures I at UMass Lowell.

Twenty-one years after the publication of the 2nd edition of *Applied Discrete Structures for Computer Science*, in 1989 the publishing and computing landscape have both changed dramatically. We signed a contract for the second edition with Science Research Associates but by the time the book was ready to print, SRA had been sold to MacMillan. Soon after, the rights had been passed on to Pearson Education, Inc. In 2010, the long-term future of printed textbooks is uncertain. In the meantime, textbook prices (both printed and e-books) have increased and a growing open source textbook market movement has started. One of our objectives in revisiting this text is to make it available to our students in an affordable format. In its original form, the text was peer-reviewed and was adopted for use at several universities throughout the country. For this reason, we see *Applied Discrete Structures* as not only an inexpensive alternative, but a high quality alternative.

As indicated above the computing landscape is very different from the 1980's and accounts for the most significant changes in the text. One of the most common programming languages of the 1980's, Pascal; and we used it to illustrate many of the concepts in the text. Although it isn't totally dead, Pascal is far from the mainstream of computing in the 21st century. In 1989, *Mathematica* had been out for less than a year — now a major force in scientific computing. The open source software movement also started in the 1980's and in 2005, the first version of Sage, an open-source alternative to *Mathematica* was first released. In *Applied Discrete Structures* we have replaced "Pascal Notes" with "Mathematica Notes" and "Sage Notes." Finally, 1989 was the year that World Wide Web was invented by Tim Berners-Lee. There wasn't a single www in the 2nd edition. In this version, we intend to make use of extensive web resources, including video demonstrations.

We would like to thank Tony Penta, Sitansu Mittra, and Dan Klain for using the preliminary versions of *Applied Discrete Structures*. The corrections and input they provided was appreciated.

We repeat the preface to *Applied Discrete Structures for Computer Science* below. Plans for the instructor's guide, which is mentioned in the preface are uncertain at this time.

Preface to Applied Discrete Structures for Computer Science, 2nd Ed.

We feel proud and fortunate that most authorities, including MAA and ACM, have settled on a discrete mathematics syllabus that is virtually identical to the contents of the first edition of *Applied Discrete Structures for Computer Science*. For that reason, very few topical changes needed to be made in this new edition, and the order of topics is almost unchanged. The main change is the addition of a large number of exercises at all levels. We have "fine-tuned" the contents by expanding the preliminary coverage of sets and combinatorics, and we have added a discussion of binary integer representation. We have also added an introduction including several examples, to provide motivation for those students who may find it reassuring to know that mathematics has "real" applications. "Appendix B—Introduction to Algorithms," has also been added to make the text more self-contained.

How This Book Will Help Students

In writing this book, care was taken to use language and examples that gradually wean students from a simpleminded mechanical approach and move them toward mathematical maturity. We also recognize that many students who hesitate to ask for help from an instructor need a readable text, and we have tried to anticipate the questions that go unasked.

The wide range of examples in the text are meant to augment the "favorite examples" that most instructors have for teaching the topics in discrete mathematics.

To provide diagnostic help and encouragement, we have included solutions and/or hints to the odd-numbered exercises. These solutions include detailed answers whenever warranted and complete proofs, not just terse outlines of proofs.

Our use of standard terminology and notation makes *Applied Discrete Structures for Computer Science* a valuable reference book for future courses. Although many advanced books have a short review of elementary topics, they cannot be complete.

How This Book Will Help Instructors

The text is divided into lecture-length sections, facilitating the organization of an instructor's presentation.

Topics are presented in such a way that students' understanding can be monitored through thought-provoking exercises. The exercises require an understanding of the topics and how they are interrelated, not just a familiarity with the key words.

An Instructor's Guide is available to any instructor who uses the text. It includes:

- (a) Chapter-by-chapter comments on subtopics that emphasize the pitfalls to avoid;
- (b) Suggested coverage times;
- (c) Detailed solutions to most even-numbered exercises;
- (d) Sample quizzes, exams, and final exams.

How This Book Will Help the Chairperson/Coordinator

The text covers the standard topics that all instructors must be aware of; therefore it is safe to adopt *Applied Discrete Structures for Computer Science* before an instructor has been selected.

The breadth of topics covered allows for flexibility that may be needed due to last-minute curriculum changes.

Since discrete mathematics is such a new course, faculty are often forced to teach the course without being completely familiar with it. An Instructor's Guide is an important feature for the new instructor.

What a Difference Five Years Makes!

In the last five years, much has taken place in regards to discrete mathematics. A review of these events is in order to see how they have affected the Second Edition of *Applied Discrete Structures for Computer Science*.

(1) Scores of discrete mathematics texts have been published. Most texts in discrete mathematics can be classified as one-semester or two-semester texts. The two-semester texts, such as *Applied Discrete Structures for Computer Science*, differ in that the logical prerequisites for a more thorough study of discrete mathematics are developed.

(2) Discrete mathematics has become more than just a computer science support course. Mathematics majors are being required to take it, often before calculus. Rather than reducing the significance of calculus, this recognizes that the material a student sees in a discrete mathematics/structures course strengthens his or her understanding of the theoretical aspects of calculus. This is particularly important for today's students, since many high school courses in geometry stress mechanics as opposed to proofs. The typical college freshman is skill-oriented and does not have a high level of mathematical maturity. Discrete mathematics is also more typical of the higher-level courses that a mathematics major is likely to take.

(3) Authorities such as MAA, ACM, and A. Ralson have all refined their ideas of what a discrete mathematics course should be. Instead of the chaos that characterized the early '80s, we now have some agreement, namely that discrete mathematics should be a course that develops mathematical maturity.

(4) Computer science enrollments have leveled off and in some cases have declined. Some attribute this to the lay-offs that have taken place in the computer industry; but the amount of higher mathematics that is needed to advance in many areas of computer science has also discouraged many. A year of discrete mathematics is an important first step in overcoming a deficiency in mathematics.

(5) The Educational Testing Service introduced its Advanced Placement Exam in Computer Science. The suggested preparation for this exam includes many discrete mathematics topics, such as trees, graphs, and recursion. This continues the trend toward offering discrete mathematics earlier in the overall curriculum.

Acknowledgments

The authors wish to thank our colleagues and students for their comments and assistance in writing and revising this text. Among those who have left their mark on this edition are Susan Assmann, Shim Berkovitz, Tony Penta, Kevin Ryan, and Richard Winslow.

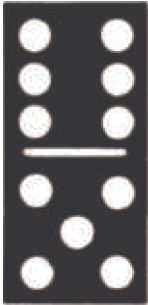
We would also like to thank Jean Hutchings, Kathy Sullivan, and Michele Walsh for work that they did in typing this edition, and our department secretaries, Mrs. Lyn Misserville and Mrs. Danielle White, whose cooperation in numerous ways has been greatly appreciated.

We are grateful for the response to the first edition from the faculty and students of over seventy-five colleges and universities. We know that our second edition will be a better learning and teaching tool as a result of their useful comments and suggestions. Our special thanks to the following reviewers: David Buchthal, University of Akron; Ronald L. Davis, Millersville University; John W. Kennedy, Pace University; Betty Mayfield, Hood College; Nancy Olmsted, Worcester State College; and Pradip Shrimani, Southern Illinois University. Finally, it has been a pleasure to work with Nancy Osman, our acquisitions editor, David Morrow, our development editor, and the entire staff at SRA.

A.W. D.

K.M.L.

chapter 11



ALGEBRAIC SYSTEMS

GOALS

The primary goal of this chapter is to make the reader aware of what an algebraic system is and how algebraic systems can be studied at different levels of abstraction. After describing the concrete, axiomatic, and universal levels, we will introduce one of the most important algebraic systems at the axiomatic level, the group. In this chapter, group theory will be a vehicle for introducing the universal concepts of isomorphism, direct product, subsystem, and generating set. These concepts can be applied to all algebraic systems. The simplicity of group theory will help the reader obtain a good intuitive understanding of these concepts. In Chapter 15, we will introduce some additional concepts and applications of group theory. We will close the chapter with a discussion of how some computer hardware and software systems use the concept of an algebraic system.

11.1 Operations

One of the first mathematical skills that we all learn is how to add a pair of positive integers. A young child soon recognizes that something is wrong if a sum has two values, particularly if his or her sum is different from the teacher's. In addition, it is unlikely that a child would consider assigning a non-positive value to the sum of two positive integers. In other words, at an early age we probably know that the sum of two positive integers is unique and belongs to the set of positive integers. This is what characterizes all binary operations on a set.

Definition: Binary Operation. Let S be a nonempty set. A binary operation on S is a rule that assigns to each ordered pair of elements of S a unique element of S . In other words, a binary operation is a function from $S \times S$ into S .

Example 11.1.1. Union and intersection are both binary operations on the power set of any universe. Addition and multiplication are binary operators on the natural numbers. Addition and multiplication are binary operations on the set of 2 by 2 real matrices, $M_{2 \times 2}(\mathbb{R})$. Division is a binary operation on some sets of numbers, such as the positive reals. But on the integers ($1/2 \notin \mathbb{Z}$) and even on the real numbers ($1/0$ is not defined), division is not a binary operation.

Notes:

- (a) We stress that the image of each ordered pair must be in S . This requirement disqualifies subtraction on the natural numbers from consideration as a binary operation, since $1 - 2$ is not a natural number. Subtraction is a binary operation on the integers.
- (b) On Notation. Despite the fact that a binary operation is a function, symbols, not letters, are used to name them. The most commonly used symbol for a binary operation is an asterisk, $*$. We will also use a diamond, \diamond , when a second symbol is needed.
- (c) If $*$ is a binary operation on S and $a, b \in S$, there are three common ways of denoting the image of the pair (a, b) . They are:

$*a\ b$	$a\ *b$	$a\ b\ *$
Prefix Form	Infix Form	Postfix Form

We are all familiar with infix form. For example, $2 + 3$ is how everyone is taught to write the sum of 2 and 3. But notice how $2 + 3$ was just described in the previous sentence! The word *sum* preceded 2 and 3. Orally, prefix form is quite natural to us. The prefix and postfix forms are superior to infix form in some respects. In Chapter 10, we saw that algebraic expressions with more than one operation didn't need parentheses if they were in prefix or postfix form. However, due to our familiarity with infix form, we will use it throughout most of the remainder of this book.

Some operations, such as negation of numbers and complementation of sets, are not binary, but unary operators.

Definition: Unary Operation. Let S be a nonempty set. A unary operator on S is a rule that assigns to each element of S a unique element of S . In other words, a unary operator is a function from S into S .

COMMON PROPERTIES OF OPERATIONS

Whenever an operation on a set is encountered, there are several properties that should immediately come to mind. To effectively make use of an operation, you should know which of these properties it has. By now, you should be familiar with most of these properties. We will list the most common ones here to refresh your memory and define them for the first time in a general setting. Let S be any set and $*$ a binary operation on S .

Properties that apply to a single binary operation:

Let $*$ be a binary operation on a set S

- $*$ is **commutative** if $a * b = b * a$ for all $a, b \in S$.
- $*$ is **associative** if $(a * b) * c = a * (b * c)$ for all $a, b, c \in S$.
- $*$ **has an identity** if there exists an element, e , in S such that $a * e = e * a = a$ for all $a \in S$.
- $*$ has the **inverse property** if for each $a \in S$, there exists $b \in S$ such that $a * b = b * a = e$.

We call b an inverse of a .

- $*$ is **idempotent** if $a * a = a$ for all $a \in S$. Properties that apply to two binary operations:

Let \diamond be a second binary operation on S .

- \diamond is **left distributive** over $*$ if $a \diamond (b * c) = (a \diamond b) * (a \diamond c)$ for all $a, b, c \in S$.
- \diamond is **right distributive** over $*$ if $(b * c) \diamond a = (b \diamond a) * (c \diamond a)$ for all $a, b, c \in S$.
- \diamond is **distributive** over $*$ if \diamond is both left and right distributive over $*$.

Let $-$ be a unary operation.

A unary operation $-$ on S has the **involution property** if $-(-a) = a$ for all $a \in S$.

Finally, a property of sets, as they relate to operations.

If T is a subset of S , we say that T is **closed** under $*$ if $a, b \in T$ implies that $a * b \in T$. In other words, by operating on elements of T with $*$, you can't obtain new elements that are outside of T .

Example 11.1.2.

- (a) The odd integers are closed under multiplication, but not under addition.
- (b) Let p be a proposition over U and let A be the set of propositions over U that imply p . That is; $q \in A$ if $q \Rightarrow p$. Then A is closed under both conjunction and disjunction.
- (c) The set positive integers that are multiples of 5 is closed under both addition and multiplication.

Note: It is important to realize that the properties listed above depend on both the set and the operation(s).

OPERATION TABLES

If the set on which an operation is defined is small, a table is often a good way of describing the operation. For example, we might want to define \oplus on $\{0, 1, 2\}$ by

$$a \oplus b = \begin{cases} a + b & \text{if } a + b < 3 \\ a + b - 3 & \text{if } a + b \geq 3 \end{cases}$$

The table for \oplus is

"

\oplus	0	1	2
0	0	1	2
1	1	2	0
2	2	0	1

The top row and left column are the column and row headings, respectively. To determine $a \oplus b$, find the entry in Row a and Column b . The following operation table serves to define $*$ on $\{i, j, k\}$.

"

$*$	i	j	k
i	i	i	i
j	j	j	j
k	k	k	k

Note that; $j * k = j$, yet $k * j = k$. Thus, $*$ is not commutative. Commutivity is easy to identify in a table: the table must be symmetric with respect to the diagonal going from the top left to lower right.

EXERCISES FOR SECTION 11.1

A Exercises

- Determine the properties that the following operations have on the positive integers.
 - addition
 - multiplication
 - M defined by $a M b = \text{larger of } a \text{ and } b$
 - m defined by $a m b = \text{smaller of } a \text{ and } b$
 - $@$ defined by $a @ b = a^b$
- Which pairs of operations in Exercise 1 are distributive over one another?
- Let $*$ be an operation on a set S and $A, B \subseteq S$. Prove that if A and B are both closed under $*$, then $A \cap B$ is also closed under $*$, but $A \cup B$ need not be.
- How can you pick out the identity of an operation from its table?
- Define $a * b$ by $|a - b|$, the absolute value of $a - b$. Which properties does $*$ have on the set of natural numbers, \mathbb{N} ?

11.2 Algebraic Systems

An algebraic system is a mathematical system consisting of a set called the domain and one or more operations on the domain. If V is the domain and $*_1, *_2, \dots, *_n$ are the operations, $[V; *_1, *_2, \dots, *_n]$ denotes the mathematical system. If the context is clear, this notation is abbreviated to V .

Example 11.2.1.

(a) Let B^* be the set of all finite strings of 0's and 1's including the null (or empty) string, λ . An algebraic system is obtained by adding the operation of concatenation. The concatenation of two strings is simply the linking of the two strings together in the order indicated. The concatenation of strings a with b is denoted $a \langle \rangle b$. For example, "01101" $\langle \rangle$ "101" = "01101101" and $\lambda \langle \rangle$ "100" = "100". Note that concatenation is an associative operation and that λ is the identity for concatenation.

Note on Notation: There isn't a standard symbol for concatenation. We have chosen $\langle \rangle$ to be consistent with the notation used in *Mathematica* for the **StringJoin** function, which does concatenation. Many programming languages use the plus sign for concatenation, but others use $\&$ or \parallel .

(b) Let M be any nonempty set and let $*$ be any operation on M that is associative and has an identity in M . Our second example might seem strange, but we include it to illustrate a point. The algebraic system $[B^*; \langle \rangle]$ is a special case of $[M; *]$. Most of us are much more comfortable with B^* than with M . No doubt, the reason is that the elements in B^* are more concrete. We know what they look like and exactly how they are combined. The description of M is so vague that we don't even know what the elements are, much less how they are combined. Why would anyone want to study M ? The reason is related to this question: What theorems are of interest in an algebraic system? Answering this question is one of our main objectives in this chapter. Certain properties of algebraic systems are called algebraic properties, and any theorem that says something about the algebraic properties of a system would be of interest. The ability to identify what is algebraic and what isn't is one of the skills that you should learn from this chapter.

Now, back to the question of why we study M . Our answer is to illustrate the usefulness of M with a theorem about M .

Theorem 11.2.1. If a, b are elements of M and $a * b = b * a$, then $(a * b) * (a * b) = (a * a) * (b * b)$.

Proof:

$$\begin{aligned} (a * b) * (a * b) &= a * (b * (a * b)) && \text{Why?} \\ &= a * ((b * a) * b) && \text{Why?} \\ &= a * ((a * b) * b) && \text{Why?} \\ &= a * (a * (b * b)) && \text{Why?} \\ &= (a * a) * (b * b) && \text{Why?} \end{aligned}$$

The power of this theorem is that it can be applied to any algebraic system that M describes. Since B^* is one such system, we can apply Theorem 11.2.1 to any two strings that commute—for example, 01 and 0101. Although a special case of this theorem could have been proven for B^* , it would not have been any easier to prove, and it would not have given us any insight into other special cases of M .

Example 11.2.2. Consider the set of 2×2 real matrices, $M_{2 \times 2}(\mathbb{R})$, with the operation of matrix multiplication. In this context, Theorem 11.2.1

can be interpreted as saying that if $AB = BA$, then $(AB)^2 = A^2 B^2$. One pair of matrices that this theorem applies to is $\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$ and

$$\begin{pmatrix} 3 & -4 \\ -4 & 3 \end{pmatrix}.$$

LEVELS OF ABSTRACTION

One of the fundamental tools in mathematics is abstraction. There are three levels of abstraction that we will identify for algebraic systems: concrete, axiomatic, and universal.

Concrete Level. Almost all of the mathematics that you have done in the past was at the concrete level. As a rule, if you can give examples of a few typical elements of the domain and describe how the operations act on them, you are describing a concrete algebraic system. Two examples of concrete systems are B^* and $M_{2 \times 2}(\mathbb{R})$. A few others are:

- (a) The integers with addition. Of course, addition isn't the only standard operation that we could include. Technically, if we were to add multiplication, we would have a different system.
- (b) The subsets of the natural numbers, with union, intersection, and complementation.
- (c) The complex numbers with addition and multiplication.

Axiomatic Level. The next level of abstraction is the axiomatic level. At this level, the elements of the domain are not specified, but certain axioms are stated about the number of operations and their properties. The system that we called M is an axiomatic system. Some combinations of axioms are so common that a name is given to any algebraic system to which they apply. Any system with the properties of M is called a *monoid*. The study of M would be called monoid theory. The assumptions that we made about M , associativity and the existence of an identity, are called the monoid axioms. One of your few brushes with the axiomatic level may have been in your elementary algebra course. Many algebra texts identify the properties of the real numbers with addition and multiplication as the field axioms. As we will see in Chapter 16, "Rings and Fields," the real numbers share these axioms with other concrete systems, all of which are called fields.

Universal Level. The final level of abstraction is the universal level. There are certain concepts, called universal algebra concepts, that can be applied to the study of all algebraic systems. Although a purely universal approach to algebra would be much too abstract for our purposes, defining concepts at this level should make it easier to organize the various algebraic theories in your own mind. In this chapter, we will consider the concepts of isomorphism, subsystem, and direct product.

GROUPS

To illustrate the axiomatic level and the universal concepts, we will consider yet another kind of axiomatic system, the group. In Chapter 5 we noted that the simplest equation in matrix algebra that we are often called upon to solve is $AX = B$, where A and B are known square matrices and X is an unknown matrix. To solve this equation, we need the associative, identity, and inverse laws. We call the systems that have these properties groups.

Definition: Group. A group consists of a nonempty set G and an operation $*$ on G satisfying the properties

- (a) $*$ is associative on G : $(a * b) * c = a * (b * c)$ for all $a, b, c \in G$.
- (b) There exists an identity element, $e \in G$ such that $a * e = e * a = a$ for all $a \in G$.
- (c) For all $a \in G$, there exists an inverse, there exist $b \in G$ such that $a * b = b * a = e$.

A group is usually denoted by its set's name, G , or occasionally by $[G; *]$ to emphasize the operation. At the concrete level, most sets have a standard operation associated with them that will form a group. As we will see below, the integers with addition is a group. Therefore, in group theory \mathbb{Z} always stands for $[\mathbb{Z}; +]$.

Generic Symbols. At the axiomatic and universal levels, there are often symbols that have a special meaning attached to them. In group theory, the letter e is used to denote the identity element of whatever group is being discussed. A little later, we will prove that the inverse of a group element, a , is unique and its inverse is usually denoted a^{-1} and is read "a inverse." When a concrete group is discussed, these symbols are dropped in favor of concrete symbols. These concrete symbols may or may not be similar to the generic symbols. For example, the identity element of the group of integers is 0, and the inverse of n is denoted by $-n$, the additive inverse of n .

The asterisk could also be considered a generic symbol since it is used to denote operations on the axiomatic level.

Example 11.2.3.

- (a) The integers with addition is a group. We know that addition is associative. Zero is the identity for addition: $0 + n = n + 0 = n$ for all integers n . The additive inverse of any integer is obtained by negating it. Thus the inverse of n is $-n$.
- (b) The integers with multiplication is not a group. Although multiplication is associative and 1 is the identity for multiplication, not all integers have a multiplicative inverse in \mathbb{Z} . For example, the multiplicative inverse of 10 is $\frac{1}{10}$, but $\frac{1}{10}$ is not an integer.
- (c) The power set of any set U with the operation of symmetric difference, \oplus , is a group. If A and B are sets, then $A \oplus B = (A \cup B) - (A \cap B)$. We will leave it to the reader to prove that \oplus is associative over $\mathcal{P}(U)$. The identity of the group is the empty set: $A \oplus \emptyset = A$. Every set is its own inverse since $A \oplus A = \emptyset$. Note that $\mathcal{P}(U)$ is not a group with union or intersection.

Definition: Abelian Group. A group is abelian if its operation is commutative.

Most of the groups that we will discuss in this book will be abelian. The term abelian is used to honor the Norwegian mathematician N. Abel (1802-29), who helped develop group theory.



Norwegian Stamp honoring Abel

EXERCISES FOR SECTION 11.2

A Exercises

- Discuss the analogy between the terms generic and concrete for algebraic systems and the terms generic and trade for prescription drugs.
- Discuss the connection between groups and monoids. Is every monoid a group? Is every group a monoid?
- Which of the following are groups?

- (a) B^* with concatenation (Example 11.2.1a).
- (b) $M_{2 \times 3}(\mathbb{R})$ with matrix addition.
- (c) $M_{2 \times 3}(\mathbb{R})$ with matrix multiplication.
- (d) The positive real numbers, \mathbb{R}^+ , with multiplication.
- (e) The nonzero real numbers, \mathbb{R}^* , with multiplication.
- (f) $\{1, -1\}$ with multiplication.
- (g) The positive integers with the operation M defined by $a M b = \text{larger of } a \text{ and } b$.

4. Prove that, \oplus , defined by $A \oplus B = (A \cup B) - (A \cap B)$ is an associative operation on $\mathcal{P}(U)$.

5. The following problem supplies an example of a non-abelian group. A rook matrix is a matrix that has only 0's and 1's as entries such that each row has exactly one 1 and each column has exactly one 1. The term rook matrix is derived from the fact that each rook matrix represents the placement of n rooks on an $n \times n$ chessboard such that none of the rooks can attack one another. A rook in chess can move only vertically or horizontally, but not diagonally. Let R_n be the set of $n \times n$ rook matrices. There are six 3×3 rook matrices:

$$I = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad R_1 = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{pmatrix} \quad R_2 = \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}$$

$$F_1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix} \quad F_2 = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix} \quad F_3 = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

- (a) List the 2×2 rook matrices. They form a group, R_2 , under matrix multiplication. Write out the multiplication table. Is the group abelian?
 - (b) Write out the multiplication table for R_3 . This is another group. Is it abelian?
 - (c) How many 4×4 rook matrices are there? How many $n \times n$ rook matrices are there?
6. For each of the following sets, identify the standard operation that results in a group. What is the identity of each group?
- (a) The set of all 2×2 matrices with real entries and nonzero determinants.
 - (b) The set of 2×3 matrices with rational entries.

B Exercises

7. Let $V = \{e, a, b, c\}$. Let $*$ be defined (partially) by $x * x = e$ for all $x \in V$. Write a complete table for $*$ so that $[V; *]$ is a group.

11.3 Some General Properties of Groups

In this section, we will present some of the most basic theorems of group theory. Keep in mind that each of these theorems tells us something about every group. We will illustrate this point at the close of the section.

Theorem 11.3.1. *The identity of a group is unique.*

One difficulty that students often encounter is how to get started in proving a theorem like this. The difficulty is certainly not in the theorem's complexity. Before actually starting the proof, we rephrase the theorem so that the implication it states is clear.

Theorem 11.3.1 (Rephrased). *If $G = [G; *]$ is a group and e is an identity of G , then no other element of G is an identity of G .*

Proof (Indirect): Suppose that $f \in G$, $f \neq e$, and f is an identity of G . We will show that $f = e$, a contradiction, which completes the proof:

$$\begin{aligned} f &= f * e \quad \text{Since } e \text{ is an identity.} \\ &= e. \quad \text{Since } f \text{ is an identity.} \quad \blacksquare \end{aligned}$$

Theorem 11.3.2. *The inverse of any element of a group is unique.*

The same problem is encountered here as in the previous theorem. We will leave it to the reader to rephrase this theorem. The proof is also left to the reader to write out in detail. Here is a hint: If b and c are both inverses of a , then you can prove that $b = c$. If you have difficulty with this proof, note that we have already proven it in a concrete setting in Chapter 5.

The significance of Theorem 11.3.2 is that we can refer to the inverse of an element without ambiguity. The notation for the inverse of a is usually a^{-1} . (note the exception below).

Example 11.3.1.

- (a) In any group, e^{-1} is the inverse of the identity e , which always is e .
- (b) $(a^{-1})^{-1}$ is the inverse of a^{-1} , which is always equal to a (see Theorem 11.3.3 below).
- (c) $(x * y * z)^{-1}$ is the inverse of $x * y * z$.
- (d) In a concrete group with an operation that is based on addition, the inverse of a is usually written $-a$. For example, the inverse of $k - 3$ in the group $[\mathbb{Z}; +]$ is written $-(k - 3) = 3 - k$. In the group of 2×2 matrices over the real numbers under matrix addition, the inverse of $\begin{pmatrix} 4 & 1 \\ 1 & -3 \end{pmatrix}$ is written $-\begin{pmatrix} 4 & 1 \\ 1 & -3 \end{pmatrix}$, which equals $\begin{pmatrix} -4 & -1 \\ -1 & 3 \end{pmatrix}$.

Theorem 11.3.3. *If a is an element of group G , then $(a^{-1})^{-1} = a$.*

Theorem 11.3.3 (Rephrased). *If a has inverse b and b has inverse c , then $a = c$.*

Proof:

$$\begin{aligned} a &= a * (b * c) \quad \text{because } c \text{ is the inverse of } b \\ &= (a * b) * c \quad \text{why?} \\ &= e * c \quad \text{why?} \\ &= c. \quad \text{by the identity property of } e. \quad \blacksquare \end{aligned}$$

Theorem 11.3.4. *If a and b are elements of group G , then $(a * b)^{-1} = b^{-1} * a^{-1}$.*

Note: This theorem simply gives you a formula for the inverse of $a * b$. This formula should be familiar. In Chapter 5 we saw that if A and B are invertible matrices, then $(AB)^{-1} = B^{-1}A^{-1}$.

Proof: Let $x = b^{-1} * a^{-1}$. We will prove that x inverts $a * b$. Since we know that the inverse is unique, we will have proved the theorem.

$$\begin{aligned} (a * b) * x &= (a * b) * (b^{-1} * a^{-1}) \\ &= a * (b * (b^{-1} * a^{-1})) \\ &= a * ((b * b^{-1}) * a^{-1}) \\ &= a * (e * a^{-1}) \\ &= a * a^{-1} \\ &= e \end{aligned}$$

Similarly, $x * (a * b) = e$; therefore, $(a * b)^{-1} = x = b^{-1} * a^{-1}$ ■

Theorem 11.3.5. Cancellation Laws. *If a , b , and c are elements of group G , both $a * b = a * c$ and $b * a = c * a$ imply that $b = c$.*

Proof: Since $a * b = a * c$, we can operate on both $a * b$ and $a * c$ on the left with a^{-1} :

$$a^{-1} * (a * b) = a^{-1} * (a * c)$$

Applying the associative property to both sides we get

$$(a^{-1} * a) * b = (a^{-1} * a) * c$$

or

$$e * b = e * c$$

and finally

$$b = c.$$

This completes the proof of the left cancellation law. The right law can be proven in exactly the same way. ■

Theorem 11.3.6. Linear Equations in a Group. If G is a group and $a, b, \in G$, the equation $a * x = b$ has a unique solution, $x = a^{-1} * b$. In addition, the equation $x * a = b$ has a unique solution, $x = b * a^{-1}$.

Proof: (for $a * x = b$):

$$\begin{aligned} a * x &= b \\ &= e * b \\ &= (a * a^{-1}) * b \\ &= a * (a^{-1} * b) \end{aligned}$$

By the cancellation law, we can conclude that $x = a^{-1} * b$.

If c and d are two solutions of the equation $a * x = b$, then $a * c = b = a * d$ and, by the cancellation law, $c = d$. This verifies that $a^{-1} * b$ is the only solution of $a * x = b$. ■

Note: Our proof of Theorem 11.3.6 was analogous to solving $4x = 9$ in the following way:

$$4x = 9 = \left(4 \cdot \frac{1}{4}\right)9 = 4\left(\frac{1}{4}9\right)$$

Therefore, by cancelling 4,

$$x = \frac{1}{4} \cdot 9 = \frac{9}{4}.$$

Exponentiation in a Group

If a is an element of a group G , then we establish the notation that

$$\begin{aligned} a * a &= a^2 \\ a * a * a &= a^3 \\ \text{etc.} \end{aligned}$$

In addition, we allow negative exponent and define, for example, $a^{-2} = (a^2)^{-1}$

Although this should be clear, proving exponentiation properties requires a more precise recursive definition:

Definition: Exponentiation in a Group. For $n \geq 0$, define a^n recursively by $a^0 = e$ and if $n > 0$, $a^n = a^{n-1} * a$. Also, if $n > 1$, $a^{-n} = (a^n)^{-1}$.

Example 11.3.2.

(a) In the group of positive real numbers with multiplication,

$$5^3 = 5^2 \cdot 5 = (5^1 \cdot 5) \cdot 5 = ((5^0 \cdot 5) \cdot 5) \cdot 5 = ((1 \cdot 5) \cdot 5) \cdot 5 = 5 \cdot 5 \cdot 5 = 125.$$

and

$$5^{-3} = (125)^{-1} = \frac{1}{125}$$

(b) In a group with addition, we use a different form of notation, reflecting the fact that in addition repeated terms are multiples, not powers. For example, in $[\mathbb{Z}; +]$, $a + a$ is written as $2a$, $a + a + a$ is written as $3a$, etc. The inverse of a multiple of a such as $-(a + a + a + a + a) = -(5a)$ is written as $(-5)a$.

Although we define, for example, $a^5 = a^4 * a$, we need to be able to extract the single factor on the left. The following lemma justifies doing precisely that.

Lemma. Let G be a group. If $b \in G$ and $n \geq 0$, then $b^{n+1} = b * b^n$, and hence $b * b^n = b^n * b$.

Proof (by induction): If $n = 0$,

$$\begin{aligned}
b^1 &= b^0 * b && \text{by the definition of exponentiation} \\
&= e * b && \text{basis for exponentiation} \\
&= b * e && \text{identity property} \\
&= b * b^0 && \text{basis for exponentiation}
\end{aligned}$$

Now assume the formula of the lemma is true for some $n \geq 0$,

$$\begin{aligned}
b^{(n+1)+1} &= b^{(n+1)} * b && \text{by the definition of exponentiation} \\
&= (b * b^n) * b && \text{by the induction hypothesis} \\
&= b * (b^n * b) && \text{associativity} \\
&= b * (b^{n+1}) && \text{definition of exponentiation} \blacksquare
\end{aligned}$$

Based on the definitions for exponentiation above, there are several properties that can be proven. They are all identical to the exponentiation properties from elementary algebra.

Theorem 11.3.7. Properties of Exponentiation. *If a is an element of a group G , and n and m are integers,*

- (a) $a^{-n} = (a^{-1})^n$ and hence $(a^n)^{-1} = (a^{-1})^n$
- (b) $a^{n+m} = a^n * a^m$
- (c) $(a^n)^m = a^{nm}$

We will leave the proofs of these properties to the interested reader. All three parts can be done by induction. For example the proof of (b) would start by defining the proposition $p(m)$, $m \geq 0$, to be $a^{n+m} = a^n * a^m$ for all n . The basis is $p(0)$: $a^{n+0} = a^n * a^0$.

Our final theorem is the only one that contains a hypothesis about the group in question. The theorem only applies to finite groups.

Theorem 11.3.8. *If G is a finite group, $|G| = n$, and a is an element of G , then there exists a positive integer m such that $a^m = e$ and $m \leq n$.*

Proof: Consider the list a, a^2, \dots, a^{n+1} . Since there are $n + 1$ elements of G in this list, there must be some duplication. Suppose that $a^p = a^q$, with $p < q$. Let $m = q - p$. Then

$$a^m = a^{q-p} = a^q * a^{-p} = a^q * (a^p)^{-1} = a^q * (a^q)^{-1} = e$$

Furthermore, since $1 \leq p < q \leq n + 1$, $m = q - p \leq n$. \blacksquare

Consider the concrete group $[\mathbb{Z}; +]$. All of the theorems that we have stated in this section except for the last one say something about \mathbb{Z} . Among the facts that we conclude from the theorems about \mathbb{Z} are:

Since the inverse of 5 is -5, the inverse of -5 is 5.

The inverse of $-6 + 71$ is $-(71) + -(-6) = -71 + 6$.

The solution of $12 + x = 22$ is $x = -12 + 22$.

$-4(6) + 2(6) = (-4 + 2)(6) = -2(6) = -(2)(6)$.

$7(4(3)) = (7 \cdot 4)(3) = 28(3)$ (twenty-eight 3s).

EXERCISES FOR SECTION 11.3

A Exercises

1. Let $[G; *]$ be a group and a be an element of G . Define $f: G \rightarrow G$ by $f(x) = a * x$.

(a) Prove that f is a bijection.

(b) On the basis of part a, describe a set of bijections on the set of integers.

2. Rephrase Theorem 11.3.2 and write out a clear proof.

3. Prove by induction on n that if a_1, a_2, \dots, a_n are elements of a group G , $n \geq 2$, then

$$(a_1 * a_2 * \dots * a_n)^{-1} = a_n^{-1} * \dots * a_2^{-1} * a_1^{-1}.$$

Interpret this result in terms of $[\mathbb{Z}; +]$ and $[\mathbb{R}; *]$.

4. True or false? If a, b, c are elements of a group G , and $a * b = c * a$, then $b = c$. Explain your answer.

5. Prove Theorem 11.3.7.

6. Each of the following facts can be derived by identifying a certain group and then applying one of the theorems of this section to it. For each fact, list the group and the theorem that are used.

- (a) $\left(\frac{1}{3}\right)5$ is the only solution of $3x = 5$.
- (b) $-(-(-18)) = -18$.
- (c) If A, B, C are 3×3 matrices over the real numbers, with $A + B = A + C$, then $B = C$.
- (d) There is only one subset of the natural numbers for which $K \oplus A = A$ for every $A \subseteq N$.

11.4 \mathbb{Z}_n , the Integers Modulo n

In this section we introduce a collection of concrete groups, one for each positive integer, that will provide us with a wealth of examples and applications. We start with a theorem about integer division that is intuitively clear. We leave the proof as an optional exercise.

The Division Property for Integers. If $m, n \in \mathbb{Z}$, $n > 0$, then there exist two unique integers, q (quotient) and r (remainder), such that $m = nq + r$ and $0 \leq r < n$.

Note: The division property says that if m is divided by n , you will obtain a quotient and a remainder, where the remainder is less than n . This is a fact that most elementary school students learn when they are introduced to long division. In doing the division problem $1986 \div 97$, you obtain a quotient of 20 and a remainder of 46. This result could either be written $\frac{1986}{97} = 20 + \frac{46}{97}$ or $1986 = 97 \cdot 20 + 46$. The later form is how the division property is normally expressed.

If two numbers, a and b , share the same remainder after dividing by n , we say that they are congruent modulo n , denoted $a \equiv b \pmod{n}$. For example, $13 \equiv 38 \pmod{5}$ because $13 = 5 \cdot 2 + 3$ and $38 = 5 \cdot 7 + 3$.

Modular Arithmetic. If n is a positive integer, we define the operations of addition modulo n ($+_n$) and multiplication modulo n (\times_n) as follows. If $a, b \in \mathbb{Z}$,

$a +_n b$ = the remainder after $a + b$ is divided by n

$a \times_n b$ = the remainder after $a \cdot b$ is divided by n .

Notes:

- The result of doing arithmetic modulo n is always an integer between 0 and $n - 1$, by the Division Property. This observation implies that $\{0, 1, \dots, n - 1\}$ is closed under modulo n arithmetic.
- It is always true that $a +_n b \equiv (a + b) \pmod{n}$ and $a \times_n b \equiv (a \cdot b) \pmod{n}$. For example, $4 +_7 5 = 2 \equiv 9 \pmod{7}$ and $4 \times_7 5 \equiv 6 \equiv 20 \pmod{7}$.
- We will use the notation \mathbb{Z}_n to denote the set $\{0, 1, 2, \dots, n - 1\}$.

Properties of Modular Arithmetic on \mathbb{Z}_n

Addition modulo n is always commutative and associative; 0 is the identity for $+_n$ and every element of \mathbb{Z}_n has an additive inverse.

Multiplication modulo n is always commutative and associative, and 1 is the identity for \times_n .

Theorem 11.4.1. If $a \in \mathbb{Z}_n$, $a \neq 0$, then the additive inverse of a is $n - a$.

Proof: $a + (n - a) = n \equiv 0 \pmod{n}$, since $n = n \cdot 1 + 0$. Therefore, $a +_n (n - a) = 0$ ■

Note: The algebraic properties of $+_n$ and \times_n on \mathbb{Z}_n are identical to the properties of addition and multiplication on \mathbb{Z} .

The Group \mathbb{Z}_n . For each $n \geq 1$, $(\mathbb{Z}_n; +_n)$ is a group. Henceforth, we will use the abbreviated notation \mathbb{Z}_n when referring to this group. Figure 11.4.1 contains the tables for \mathbb{Z}_1 through \mathbb{Z}_6 .

	0
0	0

	0	1
0	0	1
1	1	0

	0	1	2
0	0	1	2
1	1	2	0
2	2	0	1

	0	1	2	3
0	0	1	2	3
1	1	2	3	0
2	2	3	0	1
3	3	0	1	2

	0	1	2	3	4
0	0	1	2	3	4
1	1	2	3	4	0
2	2	3	4	0	1
3	3	4	0	1	2
4	4	0	1	2	3

	0	1	2	3	4	5
0	0	1	2	3	4	5
1	1	2	3	4	5	0
2	2	3	4	5	0	1
3	3	4	5	0	1	2
4	4	5	0	1	2	3
5	5	0	1	2	3	4

Figure 11.4.1
Addition tables for \mathbb{Z}_n , $1 \leq n \leq 6$.

Example 11.4.1.

- We are all somewhat familiar with \mathbb{Z}_{12} since the hours of the day are counted using this group, except for the fact that 12 is used in place of 0. Military time uses the mod 24 system and does begin at 0. If someone started a four-hour trip at hour 21, the time at which she would arrive is $21 +_{24} 4 = 1$. If a satellite orbits the earth every four hours and starts its first orbit at hour 5, it would end its first orbit at time

$5 +_{24} 4 = 9$. Its tenth orbit would end at $5 +_{24} 7 \times_{24} 4 = 9$ hours on the clock

(b) Virtually all computers represent unsigned integers in binary form with a fixed number of digits. A very small computer might reserve seven bits to store the value of an integer. There are only 2^7 different values that can be stored in seven bits. Since the smallest value is 0, represented as 0000000, the maximum value will be $2^7 - 1 = 127$, represented as 1111111. When a command is given to add two integer values, and the two values have a sum of 128 or more, overflow occurs. For example, if we try to add 56 and 95, the sum is an eight-digit binary integer 10010111. One common procedure is to retain the seven lowest-ordered digits. The result of adding 56 and 95 would be $0010111_{\text{two}} = 23 \equiv 56 + 95 \pmod{128}$. Integer arithmetic with this computer would actually be modulo 128 arithmetic.



Mathematica Note

Most computer languages have a "mod" function that computes the remainder when one integer is divided by another. *Mathematica* is no exception. To determine the remainder upon dividing 1986 by 97 we can evaluate

```
Mod[1986, 97]
```

46

A mod 6 addition function can be defined based on **Mod** with the following input:

```
Plus6[a_, b_] := Mod[a + b, 6]
```

There is a free package called *AbstractAlgebra* that is available at <http://www.central.edu/eaam/index.asp>. It contains a function that will generate the operation tables, also called *Cayley Tables*, such you see in Figure 11.4.1. First load the package, as instructed:

```
<< AbstractAlgebra`Master`
```

We can form a the group \mathbb{Z}_6 using the **FormGroupoid** function:

```
G = FormGroupoid[Range[0, 5], Plus6]
```

Groupoid({0, 1, 2, 3, 4, 5}, -Operation-)

Then the function called **CayleyTable** generates the table for the group \mathbb{Z}_6 :

```
CayleyTable[G, BodyColored -> False,
  HeadingsColored -> False, ShowExtraCayleyInformation -> False]
```

TheGroup
y

x

*	0	1	2	3	4	5
0	0	1	2	3	4	5
1	1	2	3	4	5	0
2	2	3	4	5	0	1
3	3	4	5	0	1	2
4	4	5	0	1	2	3
5	5	0	1	2	3	4

Note: The rules **BodyColored -> False**, **HeadingsColored -> False**, **ShowExtraCayleyInformation -> False** are included in the input above for easier print readability. They would not be normally included when using **CayleyTable**.

It's actually even easier to generate these tables because the family of \mathbb{Z}_n 's is part of the package. Here is the table for \mathbb{Z}_9 :

```
CayleyTable[Z[9], BodyColored → False,
  HeadingsColored → False, ShowExtraCayleyInformation → False]
```

Z[9]
y

	+	0	1	2	3	4	5	6	7	8
x	0	0	1	2	3	4	5	6	7	8
	1	1	2	3	4	5	6	7	8	0
	2	2	3	4	5	6	7	8	0	1
	3	3	4	5	6	7	8	0	1	2
	4	4	5	6	7	8	0	1	2	3
	5	5	6	7	8	0	1	2	3	4
	6	6	7	8	0	1	2	3	4	5
	7	7	8	0	1	2	3	4	5	6
	8	8	0	1	2	3	4	5	6	7



Sage Note

Sage has some extremely powerful tool for working with groups, although the operation tables of groups \mathbb{Z}_n are not all that easy to create. Here is a very simple calculation with mod 6 arithmetic.

```
R = IntegerModRing(6)
a = R(3) + R(5)*R(2)
a
1
```

There is a built in family of groups that is essentially the same as the \mathbb{Z}_n 's. Here is the one that corresponds with \mathbb{Z}_6 , where the letters a through f would be replaced with 0 through 5.

```
G=CyclicPermutationGroup(6)
G.cayley_table()
*  a b c d e f
+-----
a | a b c d e f
b | b c d e f a
c | c d e f a b
d | d e f a b c
e | e f a b c d
f | f a b c d e
```

EXERCISES FOR SECTION 11.4

A Exercises

1. Calculate:

- $7 +_8 3$
- $7 \times_8 3$
- $4 \times_8 4$
- $10 +_{12} 2$
- $6 \times_8 2 +_8 6 \times_8 5$
- $6 \times_8 (2 +_8 5)$
- $3 \times_5 3 \times_5 3 \times_5 3 \equiv 3^4 \pmod{5}$
- $2 \times_{11} 7$

- (i) $2 \times_{14} 7$
- 2. List the additive inverses of the following elements:
 - (a) 4, 6, 9 in \mathbb{Z}_{10}
 - (b) 16, 25, 40 in \mathbb{Z}_{50}
- 3. In the group \mathbb{Z}_{11} , what are:
 - (a) $3(4)$?
 - (b) $36(4)$?
 - (c) How could you efficiently compute $m(4)$, $m \in \mathbb{Z}$?
- 4. Prove that $\{1, 2, 3, 4\}$ is a group under the operation \times_5 .
- 5. A student is asked to solve the following equations under the requirement that all arithmetic should be done in \mathbb{Z}_2 . List all solutions.
 - (a) $x^2 + 1 = 0$.
 - (b) $x^2 + x + 1 = 0$.
- 6. Determine the solutions of the same equations as in Exercise 5 in \mathbb{Z}_5 .

B Exercises

- 7. Prove the division property by induction on m .
- 8. Prove that congruence modulo n is an equivalence relation on the integers. Describe the set of equivalence classes that congruence modulo n defines.

11.5 Subsystems

The subsystem is a fundamental concept of algebra at the universal level.

Definition: Subsystem. If $[V; *_1, \dots, *_n]$ is an algebraic system of a certain kind and W is a subset of V , then W is a subsystem of V if $[W; *_1, \dots, *_n]$ is an algebraic system of the same kind as V . The usual notation for " W is a subsystem of V " is $W \leq V$.

Since the definition of a subsystem is at the universal level, we can cite examples of the concept of subsystems at both the axiomatic and concrete level.

Example 11.5.1

- (a) (Axiomatic) If $[G; *]$ is a group, and H is a subset of G , then H is a subgroup of G if $[H; *]$ is a group.
- (b) (Concrete) $U = \{-1, 1\}$ is a subgroup of $[\mathbb{R}^*; \cdot]$. Take the time now to write out the multiplication table of U and convince yourself that $[U; \cdot]$ is a group.
- (c) (Concrete) The even integers, $2\mathbb{Z} = \{2k : k \text{ is an integer}\}$ is a subgroup of $[\mathbb{Z}; +]$. Convince yourself of this fact.
- (d) (Concrete) The set of nonnegative integers is not a subgroup of $[\mathbb{Z}; +]$. All of the group axioms are true for this subset except one: no positive integer has a positive additive inverse. Therefore, the inverse property is not true. Note that every group axiom must be true for a subset to be a subgroup.
- (e) (Axiomatic) If M is a monoid and P is a subset of M , then P is a submonoid of M if P is a monoid.
- (f) (Concrete) If B^* is the set of strings of 0's and 1's of length zero or more with the operation of concatenation, then two examples of submonoids of B^* are: (i) the set of strings of even length, and (ii) the set of strings that contain no 0's. The set of strings of length less than 50 is not a submonoid because it isn't closed under concatenation. Why isn't the set of strings of length 50 or more a submonoid of B^* ?

For the remainder of this section, we will concentrate on the properties of subgroups. The first order of business is to establish a systematic way of determining whether a subset of a group is a subgroup.

Theorem/Algorithm 11.5.1. To determine whether H , a subset of group $[G; *]$, is a subgroup, it is sufficient to prove:

- (a) H is closed under $*$; that is, $a, b \in H \Rightarrow a * b \in H$;
- (b) H contains the identity element for $*$; and
- (c) H contains the inverse of each of its elements; that is, $a \in H \Rightarrow a^{-1} \in H$.

Proof: Our proof consists of verifying that if the three properties above are true, then all the axioms of a group are true for $[H; *]$. By Condition (a), $*$ can be considered an operation on H . The associative, identity, and inverse properties are the axioms that are needed. The identity and inverse properties are true by Conditions (b) and (c), respectively, leaving only the associative property. Since, $[G; *]$ is a group, $a * (b * c) = (a * b) * c$ for all $a, b, c \in G$. Certainly, if this equation is true for all choices of three elements from G , it will be true for all choices of three elements from H , since H is a subset of G . ■

For every group with at least two elements, there are at least two subgroups: they are the whole group and $\{e\}$. Since these two are automatic, they are not considered very interesting and are called the improper subgroups of the group; $\{e\}$ is sometimes referred to as the trivial subgroup. All other subgroups, if there are any, are called proper subgroups.

We can apply Theorem 11.5.1 at both the concrete and axiomatic levels.

Examples 11.5.2.

- (a) (Concrete) We can verify that $2\mathbb{Z} \leq \mathbb{Z}$, as stated in Example 11.5.1. Whenever you want to discuss a subset, you must find some convenient way of describing its elements. An element of $2\mathbb{Z}$ can be described as 2 times an integer; that is, $a \in 2\mathbb{Z}$ is equivalent to $(\exists k)_{\mathbb{Z}} (a = 2k)$. Now we can verify that the three conditions of Theorem 11.5.1 are true for $2\mathbb{Z}$. First, if $a, b \in 2\mathbb{Z}$, then there exist $j, k \in \mathbb{Z}$ such that $a = 2j$ and $b = 2k$. A common error is to write something like $a = 2j$ and $b = 2j$. This would mean that $a = b$, which is not necessarily true. That is why two different variables are needed to describe a and b . Returning to our proof, we can add a and b :

$$a + b = 2j + 2k = 2(j + k).$$

Since $j + k$ is an integer, $a + b$ is an element of $2\mathbb{Z}$. Second, the identity, 0, belongs to $2\mathbb{Z}$ ($0 = 2(0)$). Finally, if $a \in 2\mathbb{Z}$ and $a = 2k$, $-a = -(2k) = 2(-k)$, and $-k \in \mathbb{Z}$, therefore, $-a \in 2\mathbb{Z}$. By Theorem 11.5.1, $2\mathbb{Z} \leq \mathbb{Z}$.

How would this argument change if you were asked to prove that $3\mathbb{Z} \leq \mathbb{Z}$? or $n\mathbb{Z} \leq \mathbb{Z}$, $n \geq 2$?

- (b) (Concrete) We can prove that $H = \{0, 3, 6, 9\}$ is a subgroup of \mathbb{Z}_{12} . First, for each ordered pair $(a, b) \in H \times H$, $a +_{12} b$ is in H . This can be checked without too much trouble since $|H \times H| = 16$. Thus we can conclude that H is closed under $+_{12}$. Second, $0 \in H$. Third, $-0 = 0$, $-3 = 9$, $-6 = 6$, and $-9 = 3$. Therefore, the inverse of each element in H is in H .

- (c) (Axiomatic) If H and K are both subgroups of a group G , then $H \cap K$ is a subgroup of G . To justify this statement, we have no concrete information to work with, only the facts that $H \leq G$ and $K \leq G$. Our proof that $H \cap K \leq G$ reflects this and is an exercise in applying the definitions of intersection and subgroup. (i) If a and b are elements of $H \cap K$, then a and b both belong to H , and since $H \leq G$, $a * b$ must be an element of H . Similarly, $a * b \in K$; therefore, $a * b \in H \cap K$. (ii) The identity of G must belong to both H and K ; hence it belongs to $H \cap K$. (iii) If $a \in H \cap K$, then $a \in H$, and since $H \leq G$, $a^{-1} \in H$. Similarly, $a^{-1} \in K$. Hence, by the theorem, $H \cap K \leq G$.

Now that this fact has been established, we can apply it to any pair of subgroups of any group. For example, since $2\mathbb{Z}$ and $3\mathbb{Z}$ are both subgroups of $[\mathbb{Z}; +]$, $2\mathbb{Z} \cap 3\mathbb{Z}$ is also a subgroup of \mathbb{Z} . Note that if $a \in 2\mathbb{Z} \cap 3\mathbb{Z}$, a must have a factor of 3; that is, there exists $k \in \mathbb{Z}$

such that $a = 3k$. In addition, a must be even, therefore k must be even. There exists $j \in \mathbb{Z}$ such that $k = 2j$, therefore $a = 3(2j) = 6j$. This shows that $2\mathbb{Z} \cap 3\mathbb{Z} \subseteq 6\mathbb{Z}$. The opposite containment can easily be established; therefore, $2\mathbb{Z} \cap 3\mathbb{Z} = 6\mathbb{Z}$.

Given a finite group, we can apply Theorem 11.3.7 to obtain a simpler condition for a subset to be a subgroup.

Theorem/Algorithm 11.5.2. If $[G; *]$ is a finite group, H is a nonempty subset of G , and you can verify that H is closed under $*$, then H is a subgroup of G .

Proof: In this proof, we demonstrate that Conditions (b) and (c) of Theorem 11.5.1 follow from the closure of H under $*$, which is Condition (a). First, select any element of H ; call it β . The powers of β : $\beta^1, \beta^2, \beta^3, \dots$ are all in H by the closure property. By Theorem 11.3.7, there exists $m, m \leq |G|$, such that $\beta^m = e$; hence $e \in H$. To prove that (c) is true, we let a be any element of H . If $a = e$, then a^{-1} is in H since $e^{-1} = e$. If $a \neq e$, $a^q = e$ for some q between 2 and $|G|$ and

$$e = a^q = a^{q-1} * a.$$

Therefore, $a^{-1} = a^{q-1}$, which belongs to H since $q - 1 \geq 1$. ■

Example 11.5.3 To determine whether $H_1 = \{0, 5, 10\}$ and $H_2 = \{0, 4, 8, 12\}$ are subgroups of \mathbb{Z}_{15} , we need only write out the addition tables (modulo 15) for these sets.

H_1
 y

+	0	5	10
0	0	5	10
5	5	10	0
10	10	0	5

H_2
 y

*	0	4	8	12
0	0	4	8	12
4	4	8	12	1
8	8	12	1	5
12	12	1	5	9

x
x

Note that H_1 is a subgroup of \mathbb{Z}_{15} . Since the interior of the addition table for H_2 contains elements that are outside of H_2 , H_2 is not a subgroup of \mathbb{Z}_{15} .

One kind of subgroup that merits special mention due to its simplicity is the cyclic subgroup.

Definition: Cyclic Subgroup Generated by an Element. If G is a group and $a \in G$, the cyclic subgroup generated by a , $\langle a \rangle$, is the set of powers of a and their inverses:

$$\langle a \rangle = \{a^n : n \in \mathbb{Z}\}$$

A subgroup H is cyclic if there exists $a \in H$ such that $H = \langle a \rangle$.

Definition: Cyclic Group. A group G is cyclic if there exists $\beta \in G$ such that $\langle \beta \rangle = G$.

Note: If the operation on G is additive, then $\langle a \rangle = \{(n)a : n \in \mathbb{Z}\}$.

Example 11.5.4.

(a) In $[\mathbb{R}; \cdot]$, $\langle 2 \rangle = \{2^n : n \in \mathbb{Z}\} = \{\dots, \frac{1}{16}, \frac{1}{8}, \frac{1}{4}, \frac{1}{2}, 1, 2, 4, 8, 16, \dots\}$.

(b) In \mathbb{Z}_{15} , $\langle 6 \rangle = \{0, 3, 6, 9, 12\}$. If G is finite, you need list only the positive powers of a up to the first occurrence of the identity to obtain all of $\langle a \rangle$. In \mathbb{Z}_{15} , the multiples of 6 are 6, $(2)6 = 12$, $(3)6 = 3$, $(4)6 = 9$, and $(5)6 = 0$. Note that $\{0, 3, 6, 9, 12\}$ is also $\langle 3 \rangle$, $\langle 9 \rangle$, and $\langle 12 \rangle$. This shows that a cyclic subgroup can have different generators.

If you want to list the cyclic subgroups of a group, the following theorem can save you some time.

Theorem 11.5.3. If a is an element of group G , then $\langle a \rangle = \langle a^{-1} \rangle$. This is an easy way of seeing that $\langle 9 \rangle$ in \mathbb{Z}_{15} equals $\langle 6 \rangle$, since $-6 = 9$.

EXERCISES FOR SECTION 11.5

A Exercises

1. Which of the following subsets of the real numbers is a subgroup of $[\mathbb{R}; +]$?

- (a) the rational numbers
- (b) the positive real numbers
- (c) $\{k/2 \mid k \text{ is an integer}\}$
- (d) $\{2^k \mid k \text{ is an integer}\}$

- (e) $\{x \mid -100 \leq x \leq 100\}$
2. Describe in simpler terms the following subgroups of \mathbb{Z} :
- (a) $5\mathbb{Z} \cap 4\mathbb{Z}$
- (b) $4\mathbb{Z} \cap 6\mathbb{Z}$ (be careful)
- (c) the only finite subgroup of \mathbb{Z}
3. Find at least two proper subgroups of R_3 , the set of 3×3 rook matrices (see Exercise 5 of Section 11.2).
4. Where should you place the following in Figure 11.5.1?
- (a) e
- (b) a^{-1}
- (c) $x * y$

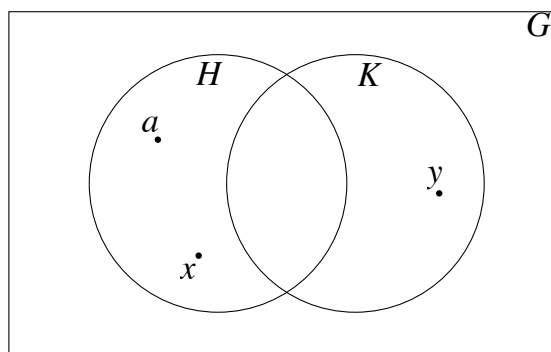


Figure 11.5.1

5. (a) List the cyclic subgroups of \mathbb{Z}_6 and draw an ordering diagram for the relation "is a subset of" on these subgroups.
- (b) Do the same for \mathbb{Z}_{12} .
- (c) Do the same for \mathbb{Z}_8 .
- (d) On the basis of your results in parts a, b, and c, what would you expect if you did the same with \mathbb{Z}_{24} ?

B Exercises

6. *Subgroups generated by subsets of a group.* The concept of a cyclic subgroup is a special case of the concept that we will discuss here. Let $[G; *]$ be a group and S a nonempty subset of G . Define the set $\langle S \rangle$ recursively by:
- (i) If $a \in S$, then $a \in \langle S \rangle$,
 - (ii) If $a, b \in \langle S \rangle$, then $a * b \in \langle S \rangle$, and
 - (iii) If $a \in \langle S \rangle$, then $a^{-1} \in \langle S \rangle$.
- (a) By its definition, $\langle S \rangle$ has all of the properties needed to be a subgroup of G . The only thing that isn't obvious is that the identity of G is in $\langle S \rangle$. Prove that the identity of G is in $\langle S \rangle$.
- (b) What is $\langle \{9, 15\} \rangle$ in $[\mathbb{Z}; +]$?
- (c) Prove that if $H \leq G$ and $S \subseteq H$, then $\langle S \rangle \leq H$. This proves that $\langle S \rangle$ is contained in every subgroup of G that contains S ; that is, $\langle S \rangle = \bigcap_{\substack{S \subseteq H \\ H \leq G}} H$.
- (d) Describe $\langle \{0.5, 3\} \rangle$ in $[\mathbb{R}^+; \cdot]$ and in $[\mathbb{R}; +]$.
- (e) If $j, k \in \mathbb{Z}$, $\langle \{j, k\} \rangle$ is a cyclic subgroup of \mathbb{Z} . In terms of j and k , what is a generator of $\langle \{j, k\} \rangle$?
7. Prove that if $H, K \leq G$, and $H \cup K = G$, then $H = G$ or $K = G$. (Hint: Use an indirect argument.)

11.6 Direct Products

Our second universal algebraic concept lets us look in the opposite direction from subsystems. Direct products allow us to create larger systems. In the following definition, we avoid complicating the notation by not specifying how many operations the systems have.

Definition: Direct Product. If $[V_1; *_1, \diamond_1, \dots], [V_2; *_2, \diamond_2, \dots], \dots, [V_n; *_n, \diamond_n, \dots]$ are algebraic systems of the same kind, then the direct product of these systems is $V = V_1 \times V_2 \times \dots \times V_n$, with operations defined below. The elements of V are n -tuples of the form (a_1, a_2, \dots, a_n) , where $a_k \in V_k, k = 1, \dots, n$. The systems V_1, V_2, \dots, V_n are called the factors of V . There are as many operations on V as there are on the factors. Each of these operations is defined componentwise:

If $(a_1, a_2, \dots, a_n), (b_1, b_2, \dots, b_n) \in V$,

$$\begin{aligned}(a_1, a_2, \dots, a_n) * (b_1, b_2, \dots, b_n) &= (a_1 *_1 b_1, a_2 *_2 b_2, \dots, a_n *_n b_n) \\ (a_1, a_2, \dots, a_n) \diamond (b_1, b_2, \dots, b_n) &= (a_1 \diamond_1 b_1, a_2 \diamond_2 b_2, \dots, a_n \diamond_n b_n) \\ &\vdots\end{aligned}$$

Example 11.6.1. Consider the monoids \mathbb{N} (the set of natural numbers with addition) and B^* (the set of finite strings of 0's and 1's with concatenation). The direct product of \mathbb{N} with B^* is a monoid. We illustrate its operation, which we will denote by $*$, with examples:

$$(4, 001) * (3, 11) = (4 + 3, 001 \text{ <> } 11) = (7, 00111)$$

$$(0, 11010) * (3, 01) = (3, 1101001)$$

$$(0, \lambda) * (129, 00011) = (0 + 129, \lambda \text{ <> } 00011) = (129, 00011)$$

$$(2, 01) * (8, 10) = (10, 0110), \text{ and}$$

$$(8, 10) * (2, 01) = (10, 1001).$$

Note that our new monoid is not commutative. What is the identity for $*$?

Notes:

(a) On notation. If two or more consecutive factors in a direct product are identical, it is common to combine them using exponential notation. For example, $\mathbb{Z} \times \mathbb{Z} \times \mathbb{R}$ can be written $\mathbb{Z}^2 \times \mathbb{R}$, and $\mathbb{R} \times \mathbb{R} \times \mathbb{R} \times \mathbb{R}$ can be written \mathbb{R}^4 . This is purely a notational convenience; no exponentiation is really taking place.

(b) In our definition of a direct product, the operations are called componentwise operations, and they are indeed operations on V . Consider $*$ above. If two n -tuples, a and b , are selected from V , the first components of a and b , a_1 and b_1 , are operated on with $*_1$ to obtain $a_1 *_1 b_1$, the first component of $a * b$. Note that since $*_1$ is an operation on V_1 , $a_1 *_1 b_1$ is an element of V_1 . Similarly, all other components of $a * b$, as they are defined, belong to their proper sets.

One significant fact about componentwise operations is that the components of the result can all be computed at the same time (concurrently). The time required to compute in a direct product can be reduced to a length of time that is not much longer than the maximum amount of time needed to compute in the factors (see Figure 11.6.1).

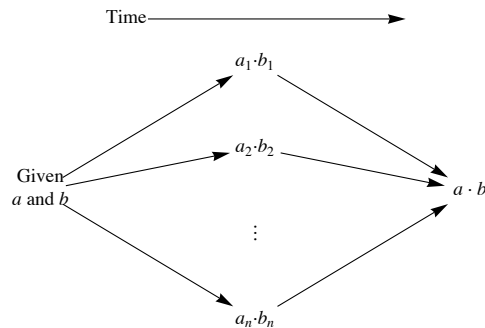


Figure 11.6.1
Concurrent calculation in a direct product.

(c) A direct product of algebraic systems is not always an algebraic system of the same type as its factors. This is due to the fact that certain axioms that are true for the factors may not be true for the set of n -tuples. This situation does not occur with groups however. You will find that whenever a new type of algebraic system is introduced, call it type T , one of the first theorems that is usually proven, if possible, is that the direct product of two or more systems of type T is a system of type T .

Theorem 11.6.1. The direct product of two or more groups is a group; that is, the algebraic properties of a system obtained by taking the direct product of two or more groups includes the group axioms.

We will only present the proof of this theorem for the direct product of two groups. Some slight revisions can be made to obtain a proof for any number of factors.

Proof: Stating that the direct product of two groups is a group is a short way of saying that if $[G_1; *_1]$ and $[G_2; *_2]$ are groups, then $[G_1 \times G_2; *]$ is also a group, where $*$ is the componentwise operation on $G_1 \times G_2$.

Associativity of $*$: If $a, b, c \in G_1 \times G_2$,

$$\begin{aligned}
 a * (b * c) &= (a_1, a_2) * ((b_1, b_2) * (c_1, c_2)) \\
 &= (a_1, a_2) * (b_1 * c_1, b_2 * c_2) \\
 &= (a_1 * (b_1 * c_1), a_2 * (b_2 * c_2)) \\
 &= ((a_1 * b_1) * c_1, (a_2 * b_2) * c_2) \\
 &= (a_1 * b_1, a_2 * b_2) * (c_1, c_2) \\
 &= ((a_1, a_2) * (b_1, b_2)) * (c_1, c_2) \\
 &= (a * b) * c
 \end{aligned}$$

Notice how the associativity property hinges on the associativity in each factor.

An identity for $*$: As you might expect, if e_1 and e_2 are identities for G_1 and G_2 , respectively, then $e = (e_1, e_2)$ is the identity for $G_1 \times G_2$. If $a \in G_1 \times G_2$,

$$\begin{aligned}
 a * e &= (a_1, a_2) * (e_1, e_2) \\
 &= (a_1 * e_1, a_2 * e_2) \\
 &= (a_1, a_2) \\
 &= a
 \end{aligned}$$

Similarly, $e * a = a$.

Inverses in $G_1 \times G_2$: The inverse of an element is determined componentwise $a^{-1} = (a_1, a_2)^{-1} = (a_1^{-1}, a_2^{-1})$. To verify, we compute $a * a^{-1}$:

$$\begin{aligned}
 a * a^{-1} &= (a_1, a_2) * (a_1^{-1}, a_2^{-1}) \\
 &= (a_1 * a_1^{-1}, a_2 * a_2^{-1}) \\
 &= (e_1, e_2) \\
 &= e
 \end{aligned}$$

Similarly, $a^{-1} * a = e$. ■

Example 11.6.2.

(a) If $n \geq 2$, \mathbb{Z}_2^n , the direct product of n factors of \mathbb{Z}_2 , is a group with 2^n elements. We will take a closer look at $\mathbb{Z}_2^3 = \mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_2$. The elements of this group are triples of zeros and ones. Since the operation on \mathbb{Z}_2 is $+$, we will use the symbol $+$ for the operation on \mathbb{Z}_2^3 . Two of the eight triples in the group are $a = (1, 0, 1)$ and $b = (0, 0, 1)$. Their "sum" is $a + b = (1 +_2 0, 0 +_2 0, 1 +_2 1) = (1, 0, 0)$. One interesting fact about this group is that each element is its own inverse. For example $a + a = (1, 0, 1) + (1, 0, 1) = (0, 0, 0)$; therefore $-a = a$. We use the additive notation for the inverse of a because we are using a form of addition. Note that $\{(0, 0, 0), (1, 0, 1)\}$ is a subgroup of \mathbb{Z}_2^3 . Write out the "addition" table for this set and apply Theorem 11.5.2. The same can be said for any set consisting of $(0, 0, 0)$ and another element of \mathbb{Z}_2^3 .

(b) The direct product of the positive real numbers with the integers modulo 4, $\mathbb{R}^+ \times \mathbb{Z}_4$ is an infinite group since one of its factors is infinite. The operations on the factors are multiplication and modular addition, so we will select the neutral symbol \diamond for the operation on $\mathbb{R}^+ \times \mathbb{Z}_4$. If $a = (4, 3)$ and $b = (0.5, 2)$, then

$$a \diamond b = (4, 3) \diamond (0.5, 2) = (4 \cdot 0.5, 3 +_4 2) = (2, 1)$$

$$b^2 = b \diamond b = (0.5, 2) \diamond (0.5, 2) = (0.25, 0),$$

$$a^{-1} = (4^{-1}, -3) = (0.25, 1) \text{ and}$$

$$b^{-1} = (0.5^{-1}, -2) = (2, 2).$$

It would be incorrect to say that \mathbb{Z}_4 is a subgroup of $\mathbb{R}^+ \times \mathbb{Z}_4$, but there is a subgroup of the direct product that closely resembles \mathbb{Z}_4 . It is $\{(1, 0), (1, 1), (1, 2), (1, 3)\}$. Its table is

\diamond	$\{1, 0\}$	$\{1, 1\}$	$\{1, 2\}$	$\{1, 3\}$
$\{1, 0\}$	$\{1, 0\}$	$\{1, 1\}$	$\{1, 2\}$	$\{1, 3\}$
$\{1, 1\}$	$\{1, 1\}$	$\{1, 2\}$	$\{1, 3\}$	$\{1, 0\}$
$\{1, 2\}$	$\{1, 2\}$	$\{1, 3\}$	$\{1, 0\}$	$\{1, 1\}$
$\{1, 3\}$	$\{1, 3\}$	$\{1, 0\}$	$\{1, 1\}$	$\{1, 2\}$

Imagine erasing $(1,)$ throughout the table and writing $+_4$ in place of \diamond . What would you get? We will explore this phenomenon in detail in the next section.

The whole direct product could be visualized as four parallel half-lines labeled 0, 1, 2, and 3 (Figure 11.6.2). On the k th line, the point that lies x units to the right of the zero mark would be (x, k) . The set $\{(2^n, (n) 1) \mid n \in \mathbb{Z}\}$, which is plotted on Figure 11.6.2, is a subgroup of $\mathbb{R}^+ \times \mathbb{Z}_4$. What cyclic subgroup is it?

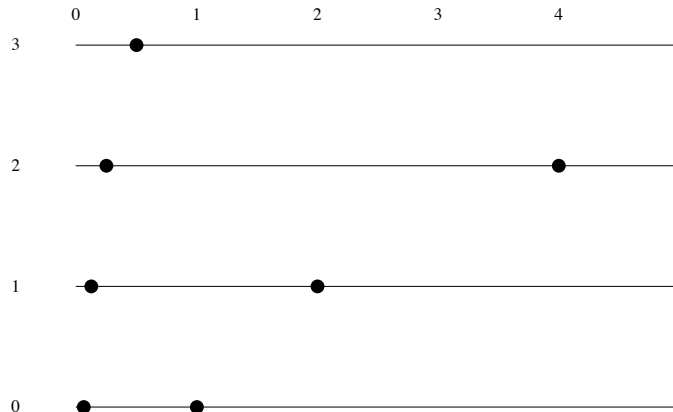


Figure 11.6.2
Graph of $\mathbb{R}^+ \times \mathbb{Z}_4$

The answer: $((2, 1))$ or $((j, 3))$.

A more conventional direct product is \mathbb{R}^2 , the direct product of two factors of $[\mathbb{R}; +]$. The operation on \mathbb{R}^2 is componentwise addition; hence we will use $+$ as the operation symbol for this group. You should be familiar with this operation, since it is identical to addition of 2×1 matrices. The Cartesian coordinate system can be used to visualize \mathbb{R}^2 geometrically. We plot the pair (s, t) on the plane in the usual way: s units along the x axis and t units along the y axis. There is a variety of different subgroups of \mathbb{R}^2 , a few of which are:

- (1) $\{(x, 0) \mid x \in \mathbb{R}\}$, all of the points on the x axis;
- (2) $\{(x, y) \mid 2x - y = 0\}$, all of the points that are on the line $2x - y = 0$;
- (3) If $a, b \in \mathbb{R}$, $\{(x, y) \mid ax + by = 0\}$. The first two subgroups are special cases of this one, which represents any line that passes through the origin.
- (4) $\{(x, y) \mid 2x - y = k, k \in \mathbb{Z}\}$, a set of lines that are parallel to $2x - y = 0$.
- (5) $\{(n, 3n) \mid n \in \mathbb{Z}\}$, which is the only countable subgroup that we have listed.

We will leave it to the reader to verify that these sets are subgroups. We will only point out how the fourth example, call it H , is closed under "addition." If $a = (p, q)$ and $b = (s, t)$ and both belong to H , then $2p - q = j$ and $2s - t = k$, where both j and k are integers.

$$a + b = (p, q) + (s, t) = (p + s, q + t)$$

We can determine whether $a + b$ belongs to H by deciding whether or not $2(p + s) - (q + t)$ is an integer:

$$\begin{aligned} 2(p + s) - (q + t) &= 2p + 2s - q - t \\ &= (2p - q) + (2s - t) \\ &= j + k \end{aligned}$$

which is an integer. This completes a proof that H is closed under the operation of \mathbb{R}^2 .

Several useful facts can be stated in regards to the direct product of two or more groups. We will combine them into one theorem, which we will present with no proof. Parts a and c were derived for $n = 2$ in the proof of Theorem 11.6.1.

Theorem 11.6.2. If $G = G_1 \times G_2 \times \cdots \times G_n$ is a direct product of n groups and $(a_1, a_2, \dots, a_n) \in G$, then:

- (a) The identity of G is (e_1, e_2, \dots, e_n) , where e_k is the identity of G_k .
- (b) $(a_1, a_2, \dots, a_n)^{-1} = (a_1^{-1}, a_2^{-1}, \dots, a_n^{-1})$.
- (c) $(a_1, a_2, \dots, a_n)^m = (a_1^m, a_2^m, \dots, a_n^m)$ for all $m \in \mathbb{Z}$.
- (d) G is abelian if and only if each of the factors G_1, G_2, \dots, G_n is abelian.
- (e) If H_1, H_2, \dots, H_n are subgroups of the corresponding factors, then $H_1 \times H_2 \times \cdots \times H_n$ is a subgroup of G .

Not all subgroups of a direct product are obtained as in part e of Theorem 11.6.2. For example, $\{(n, n) \mid n \in \mathbb{Z}\}$ is a subgroup of \mathbb{Z}^2 , but is not a direct product of two subgroups of \mathbb{Z} .

Example 11.6.3. Using the identity $(x + y) + x = y$, in \mathbb{Z}_2 , we can devise a scheme for representing a symmetrically linked list using only one link field. A symmetrically linked list is a list in which each node contains a pointer to its immediate successor and its immediate predecessor (see Figure 11.6.3). If the pointers are n -digit binary addresses, then each pointer can be taken as an element of \mathbb{Z}_2^n . Lists of this type can be accomplished using cells with only one link. In place of a left and a right pointer, the only "link" is the value of the sum (left link) + (right link). All standard list operations (merge, insert, delete, traverse, and so on) are possible with this structure, provided that you know the value of the nil pointer and the address, f , of the first (i. e., leftmost) cell. Since first f .left is nil, we can recover f .right by adding the value of nil:

$f + \text{nil} = (\text{nil} + f.\text{right}) + \text{nil} = f.\text{right}$, which is the address of the second item. Now if we temporarily retain the address, s , of the second cell, we can recover the address of the third item. The link field of the second item contains the sum $s.\text{left} + s.\text{right} = \text{first} + \text{third}$. Therefore

$$\begin{aligned} (\text{first} + \text{third}) + \text{first} &= s + s.\text{left} \\ &= (s.\text{left} + s.\text{right}) + s.\text{left} \\ &= s.\text{right} = \text{third} \end{aligned}$$

We no longer need the address of the first cell, only the second and third, to recover the fourth address, and so forth.

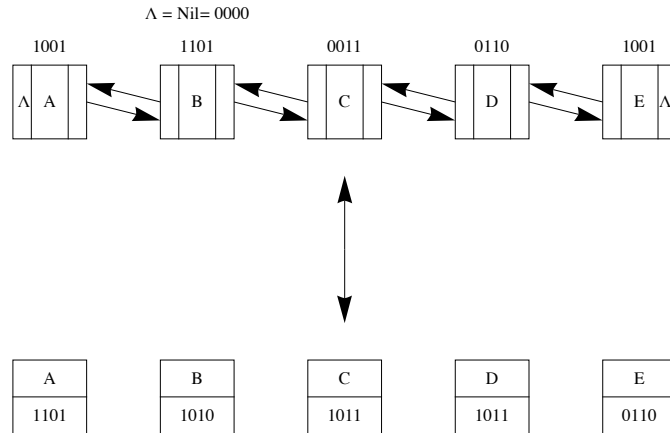


Figure 11.6.3
Symmetric Linked List

The following more formal algorithm uses names that the timing of the visits.

Algorithm 11.6.1. Given a symmetric list represented as in Example 11.6.3, a traversal of the list is accomplished as follows, where *first* is the address of the first cell. We presume that each item has some information that is represented by *item.info* and a field called *item.link* that is the sum of the left and right links.

- (1) *yesterday* = nil
- (2) *today* = *first*
- (3) While *today* \neq nil do
 - (3.1) Write(*today.info*)
 - (3.2) *tomorrow* = *today.link* + *yesterday*
 - (3.3) *yesterday* = *today*
 - (3.4) *today* = *tomorrow*.

At any point in this algorithm it would be quite easy to insert a cell between *today* and *tomorrow*. Can you describe how this would be accomplished?

EXERCISES FOR SECTION 11.6

A Exercises

1. Write out the group table of $\mathbb{Z}_2 \times \mathbb{Z}_3$ and find the two proper subgroups of this group.
2. List more examples of proper subgroups of \mathbb{R}^2 that are different from the ones in Example 11.6.2.
3. Algebraic properties of the n -cube:
 - (a) The four elements of \mathbb{Z}_2^2 can be visualized geometrically as the four corners of the 2-cube (see Figure 9.4.5). Algebraically describe the statements:
 - (i) Corners a and b are adjacent.
 - (ii) Corners a and b are diagonally opposite one another.
 - (b) The eight elements of \mathbb{Z}_2^3 can be visualized as the eight corners of the 3-cube. One face contains $\mathbb{Z}_2 \times \mathbb{Z}_2 \times \{0\}$ and the opposite face contains the remaining four elements so that $(a, b, 1)$ is behind $(a, b, 0)$. As in part a, describe statements i and ii algebraically.
 - (c) If you could imagine a geometric figure similar to the square or cube in n dimensions, and its corners were labeled by elements of \mathbb{Z}_2^n as in parts a and b, how would statements i and ii be expressed algebraically?
4. (a) Suppose that you were to be given a group $[G; *]$ and asked to solve the equation $x * x = e$. Without knowing the group, can you anticipate how many solutions there will be?
 - (b) Answer the same question as part a for the equation $x * x = x$.
5. Which of the following sets are subgroups of $\mathbb{Z} \times \mathbb{Z}$? Give a reason for any negative answers.

- (a) $\{0\}$
 - (b) $\{(2j, 2k) \mid j, k \in \mathbb{Z}\}$
 - (c) $\{(2j+1, 2k) \mid j, k \in \mathbb{Z}\}$
 - (d) $\{(n, n^2) \mid n \in \mathbb{Z}\}$
 - (e) $\{(j, k) \mid j+k \text{ is even}\}$
6. Determine the following values in group $\mathbb{Z}_3 \times \mathbb{R}^*$:
- (a) $(2, 1) * (1, 2)$
 - (b) the identity element
 - (c) $(1, 1/2)^{-1}$

1.7 Isomorphisms

The following informal definition of isomorphic systems should be memorized. No matter how technical a discussion about isomorphic systems becomes, keep in mind that this is the essence of the concept.

Definition: Isomorphic Systems/Isomorphism. Two algebraic systems are isomorphic if there exists a translation rule between them so that any true statement in one system can be translated to a true statement in the other

Example 11.7.1. Imagine that you are an eight-year-old child who has been reared in an English-speaking family, has moved to Greece, and has been placed in a Greek school. Suppose that your new teacher asks the class to do the following addition problem that has been written out in Greek.

τρία συν τέσσερα ισούται ____

The natural thing for you to do is to take out your Greek-English/English-Greek dictionary and translate the Greek words to English, as outlined in Figure 11.7.1. After you've solved the problem, you can consult the same dictionary to obtain the proper Greek word that the teacher wants. Although this is not the recommended method of learning a foreign language, it will surely yield the correct answer to the problem. Mathematically, we may say that the system of Greek integers with addition (συν) is isomorphic to English integers with addition (plus). The problem of translation between natural languages is more difficult than this though, because two complete natural languages are not isomorphic, or at least the isomorphism between them is not contained in a simple dictionary.

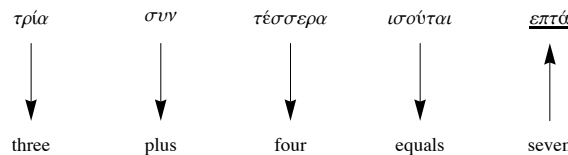


Figure 11.7.1
Solution of a Greek arithmetic problem

Example 11.7.2. Software Implementation of Sets. In this example, we will describe how set variables can be implemented on a computer. We will describe the two systems first and then describe the isomorphism between them.

System 1: The power set of $\{1, 2, 3, 4, 5\}$ with the operation union, \cup . For simplicity, we will only discuss union. However, the other operations are implemented in a similar way.

System 2: Strings of five bits of computer memory with an OR gate. Individual bit values are either zero or one, so the elements of this system can be visualized as sequences of five 0's and 1's. An OR gate, Figure 11.7.2, is a small piece of computer hardware that accepts two bit values at any one time and outputs either a zero or one, depending on the inputs. The output of an OR gate is one, except when the two bit values that it accepts are both zero, in which case the output is zero. The operation on this system actually consists of sequentially inputting the values of two bit strings into the OR gate. The result will be a new string of five 0's and 1's. An alternate method of operating in this system is to use five OR gates and to input corresponding pairs of bits from the input strings into the gates concurrently.

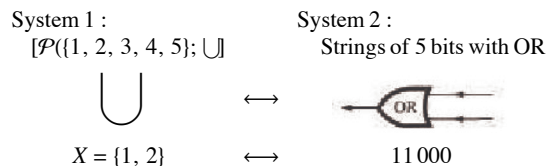


Figure 11.7.2
Translation between sets and strings of bits

The Isomorphism: Since each system has only one operation, it is clear that union and the OR gate translate into one another. The translation between sets and bit strings is easiest to describe by showing how to construct a set from a bit string. If $a_1 a_2 a_3 a_4 a_5$, is a bit string in System 2, the set that it translates to contains the number k if and only if a_k equals 1. For example, 10001 is translated to the set $\{1, 5\}$, while the set $\{1, 2\}$ is translated to 11000. Now imagine that your computer is like the child who knows English and must do a Greek problem. To execute a program that has code that includes the set expression $\{1, 2\} \cup \{1, 5\}$, it will follow the same procedure as the child to obtain the result, as shown in Figure 11.7.3.

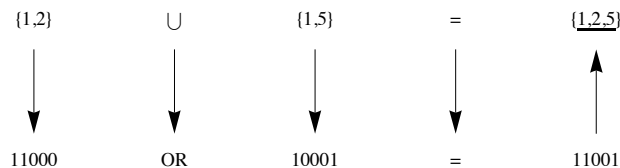


Figure 11.7.3
Translation of a problem in set theory

Example 11.7.3. Multiplying without doing multiplication. This isomorphism is between $[\mathbb{R}^+; \cdot]$ and $[\mathbb{R}; +]$. Until the 1970s, when the price of calculators dropped, multiplication and exponentiation were performed with an isomorphism between these systems. The isomorphism $(\mathbb{R}^+ \rightarrow \mathbb{R})$ between the two groups is that \cdot is translated into $+$ and any positive real number a is translated to the logarithm of a . To translate back from \mathbb{R} to \mathbb{R}^+ , you invert the logarithm function. If base ten logarithms are used, an element of \mathbb{R} , b , will be translated to 10^b . In pre-calculator days, the translation was done with a table of logarithms or with a slide rule. An example of how the isomorphism is used appears in Figure 11.7.4.

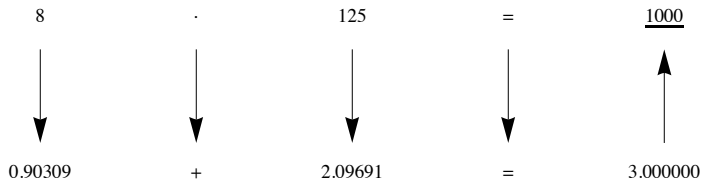


Figure 11.7.4
Multiplication using logarithms

The following definition of an isomorphism between two groups is a more formal one that appears in most abstract algebra texts. At first glance, it appears different, it is really a slight variation on the informal definition. It is the common definition because it is easy to apply; that is, given a function, this definition tells you what to do to determine whether that function is an isomorphism.

Procedure for showing that two groups are isomorphic

Definition: Group Isomorphism. If $[G_1; *_1]$ and $[G_2; *_2]$ are groups, $f: G_1 \rightarrow G_2$ is an isomorphism from G_1 into G_2 if:

- (a) f is a bijection, and
- (b) $f(a *_1 b) = f(a) *_2 f(b)$ for all $a, b \in G_1$

If such a function exists, then G_1 is isomorphic to G_2 .

Notes:

- (a) There could be several different isomorphisms between the same pair of groups. Thus, if you are asked to demonstrate that two groups are isomorphic, your answer need not be unique.
- (b) Any application of this definition requires a procedure outlined in Figure 11.7.5.

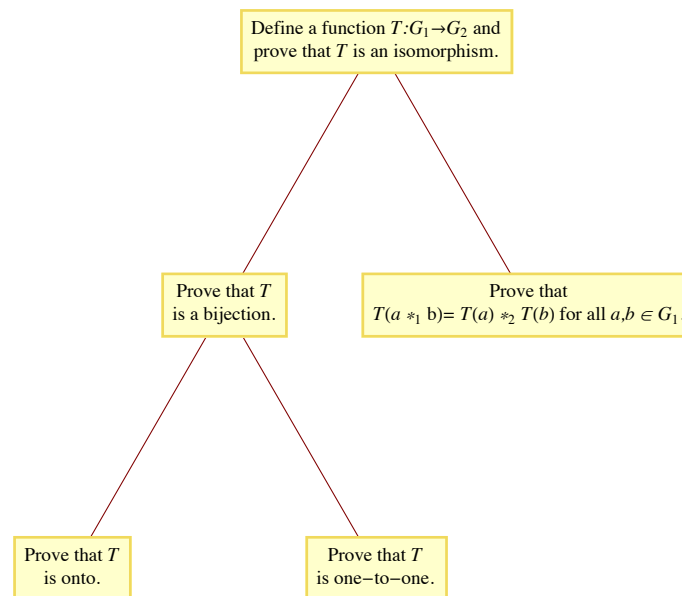


Figure 11.7.5
Steps in proving that G_1 and G_2 are isomorphic

The first condition, that an isomorphism be a bijection, reflects the fact that every true statement in the first group should have exactly one corresponding true statement in the second group. This is exactly why we run into difficulty in translating between two natural languages. To see how Condition (b) of the formal definition is consistent with the informal definition, consider the Function $L: \mathbb{R}^+ \rightarrow \mathbb{R}$ defined by $L(x) = \log_{10} x$. The translation diagram between \mathbb{R}^+ and \mathbb{R} for the multiplication problem $a \cdot b$ appears in Figure 11.7.6. We arrive at the same result by computing $L^{-1}(L(a) + L(b))$ as we do by computing $a \cdot b$. If we apply the function L to the two results, we get the same image:

$$L(a \cdot b) = L(L^{-1}(L(a) + L(b))) = L(a) + L(b) \quad (11.7a)$$

since $L(L^{-1}(x)) = x$. Note that 11.7a is exactly Condition b of the formal definition applied to the two groups \mathbb{R}^+ and \mathbb{R} .

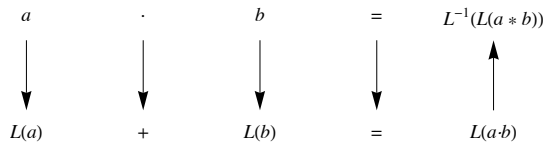


Figure 11.7.6
Multiplication using logarithms - general situation

Example 11.7.4. Consider $G = \left\{ \begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix} \mid a \in \mathbb{R} \right\}$ with matrix multiplication. This group $[\mathbb{R}; +]$ is isomorphic to G . Our translation rule is the function $f: \mathbb{R} \rightarrow G$ defined by $f(a) = \begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix}$. Since groups have only one operation, there is no need to state explicitly that addition is translated to matrix multiplication. That f is a bijection is clear from its definition. If a and b are any real numbers,

$$\begin{aligned} f(a)f(b) &= \begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & b \\ 0 & 1 \end{pmatrix} \\ &= \begin{pmatrix} 1 & a+b \\ 0 & 1 \end{pmatrix} \\ &= f(a+b) \end{aligned}$$

We can apply this translation rule to determine the inverse of a matrix in G . We know that $a + (-a) = 0$ is a true statement in \mathbb{R} . Using f to translate this statement, we get

$$\begin{aligned} f(a)f(-a) &= f(0) \\ \text{or} \\ \begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & -a \\ 0 & 1 \end{pmatrix} &= \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \end{aligned}$$

therefore,

$$\begin{pmatrix} 1 & a \\ 0 & 1 \end{pmatrix}^{-1} = \begin{pmatrix} 1 & -a \\ 0 & 1 \end{pmatrix}$$

Theorem 11.7.1 summarizes some of the general facts about group isomorphisms that are used most often in applications. We leave the proof to the reader.

Theorem 11.7.1. If $[G; *]$ and $[H, \diamond]$ are groups with identities e and e' , respectively, and $T: G \rightarrow H$ is an isomorphism from G into H , then:

- (a) $T(e) = e'$,
- (b) $T(a)^{-1} = T(a^{-1})$ for all $a \in G$, and
- (c) If K is a subgroup of G , then $T(K) = \{T(a) : a \in K\}$ is a subgroup of H and is isomorphic to K .

"Is isomorphic to" is an equivalence relation on the set of all groups. Therefore, the set of all groups is partitioned into equivalence classes, each equivalence class containing groups that are isomorphic to one another.

Procedures for showing groups are not isomorphic

How do you decide that two groups are *not* isomorphic to one another? The negation of " G and H are isomorphic" is that no translation rule between G and H exists. If G and H have different cardinalities, then no bijection from G into H can exist. Hence they are not isomorphic. Given that $|G| = |H|$, it is usually impractical to list all bijections from G into H and show that none of them satisfy Condition b of the formal definition. The best way to prove that two groups are not isomorphic is to find a true statement about one group that is not true about the other group. We illustrate this method in the following checklist that you can apply to most pairs of non-isomorphic groups in this book.

Assume that $[G; *]$ and $[H; \diamond]$ are groups. The following are reasons for G and H to be not isomorphic.

- (a) G and H do not have the same cardinality. For example, $\mathbb{Z}_{12} \times \mathbb{Z}_5$ can't be isomorphic to \mathbb{Z}_{50} and $[\mathbb{R}; +]$ can't be isomorphic to $[\mathbb{Q}^+; \cdot]$.
- (b) G is abelian and H is not abelian since $a * b = b * a$ is always true in G , but $T(a) \diamond T(b) = T(b) \diamond T(a)$ would not always be true. Two groups with six elements each are \mathbb{Z}_6 and the set of 3×3 rook matrices (see Exercise 5 in Section 11.2). The second group is non-abelian, therefore it can't be isomorphic to \mathbb{Z}_6 .
- (c) G has a certain kind of subgroup that H doesn't have. Theorem 11.7.1(c) states that this cannot happen if G is isomorphic to H . $[\mathbb{R}^+; \cdot]$ and $[\mathbb{R}^+; +]$ are not isomorphic since \mathbb{R}^+ has a subgroup with two elements, $\{-1, 1\}$, while the proper subgroups of \mathbb{R}^+ are all infinite (Convince yourself of this fact!).
- (d) The number of solutions of $x * x = e$ in G is not equal to the number of solutions of $y \diamond y = e'$ in H . \mathbb{Z}_8 is not isomorphic to \mathbb{Z}_2^3 since $x +_8 x = 0$ has two solutions, 0 and 4, while $y + y = (0, 0, 0)$ is true for all $y \in \mathbb{Z}_2^3$. If the operation in G is defined by a table, then the number of solutions of $x * x = e$ will be the number of occurrences of e in the main diagonal of the table. The equations $x^3 = e$, $x^4 = e$, ... can also be used in the same way to identify non-isomorphic groups.
- (e) One of the cyclic subgroups of G equals G (i. e., G is cyclic), while none of H 's cyclic subgroups equals H (i. e., H is noncyclic). This is a special case of Condition c. \mathbb{Z} and $\mathbb{Z} \times \mathbb{Z}$ are not isomorphic since $\mathbb{Z} = \langle 1 \rangle$ and $\mathbb{Z} \times \mathbb{Z}$ is not cyclic.

EXERCISES FOR SECTION 11.7**A Exercises**

1. State whether each pair of groups below is isomorphic. If it is, give an isomorphism; if it is not, give your reason.

(a) $\mathbb{Z} \times \mathbb{R}$ and $\mathbb{R} \times \mathbb{Z}$

(b) $\mathbb{Z}_2 \times \mathbb{Z}$ and $\mathbb{Z} \times \mathbb{Z}$

(c) \mathbb{R} and $\mathbb{Q} \times \mathbb{Q}$

(d) $\mathcal{P}(\{1, 2\})$ with symmetric difference and \mathbb{Z}_2^2

(e) \mathbb{Z}_2^2 and \mathbb{Z}_4

(f) \mathbb{R}^4 and $M_{2 \times 2}(\mathbb{R})$ with matrix addition

(g) \mathbb{R}^2 and $\mathbb{R} \times \mathbb{R}^+$

(h) \mathbb{Z}_2 and the 2×2 rook matrices

(i) \mathbb{Z}_6 and $\mathbb{Z}_2 \times \mathbb{Z}_3$

2. If you know two natural languages, show that they are not isomorphic.

3. Prove that the relation "is isomorphic to" on groups is transitive.

4. (a) Write out the operation table for $G = [\{1, -1, i, -i\}, \cdot]$ where i is the complex number for which $i^2 = -1$. Show that G is isomorphic to $[\mathbb{Z}_4; +_4]$.

(b) Solve $x^2 = -1$ in G by first translating to \mathbb{Z}_4 , solving the equation in \mathbb{Z}_4 , and then translating back to G .

B Exercises

5. It can be shown that there are five non-isomorphic groups of order eight. You should be able to describe at least three of them. Do so without use of tables. Be sure to explain why they are not isomorphic.

6. Prove Theorem 11.7.1.

7. Prove that all infinite cyclic groups are isomorphic to \mathbb{Z} .

8. (a) Prove that \mathbb{R}^* is isomorphic to $\mathbb{Z}_2 \times \mathbb{R}$.

(b) Describe how multiplication of nonzero real numbers can be accomplished doing only additions and translations.

9. Prove that if G is any group and g is some fixed element of G , then the function ϕ_g defined by $\phi_g(x) = g * x * g^{-1}$ is an isomorphism from G into itself. An isomorphism of this type is called an *automorphism*.

11.8 Using Computers to Study Groups

Groups in *Mathematica*

Mathematica has a wide variety of computable databases available and one of them is on finite groups. To access the database you use the function **FiniteGroupData**. Extensive documentation is available at <http://www.wolfram.com/mathematica/doc/mathematica/finitegroups.html>. Since we've only scratch the surface of group theory at this point, most of the groups and concepts mentioned are likely to be unfamiliar to the reader. For this reason, we will wait until Chapter 15 to discuss that database.

The *Combinatorica* package that is included in all *Mathematica* distributions has limited abstract algebra

```
<< Combinatorica`
```

Here is how to generate the body of the operation table for the ring $[Z_7; +_7]$. Notice that this really an addition table even though the function that creates the table is called **MultiplicationTable**.

```
MultiplicationTable[Range[0, 6], Function[{a, b}, Mod[a + b, 7]]]
```

$$\begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ 2 & 3 & 4 & 5 & 6 & 7 & 1 \\ 3 & 4 & 5 & 6 & 7 & 1 & 2 \\ 4 & 5 & 6 & 7 & 1 & 2 & 3 \\ 5 & 6 & 7 & 1 & 2 & 3 & 4 \\ 6 & 7 & 1 & 2 & 3 & 4 & 5 \\ 7 & 1 & 2 & 3 & 4 & 5 & 6 \end{pmatrix}$$

An even more user-friendly package that you would need to download to use is available at Exploring Abstract Algebra with Mathematica (<http://www.central.edu/EAAM/>). The package, when installed on your computer, is loaded with the command

```
<< AbstractAlgebra`Master`
```

The group Z_{12} is

```
G = ZG[12]
```

```
Groupoid({0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11}, (#1 + #2) mod 12 &)
```

At this point **G** is an object that consists of the set $\{0, 1, 2, \dots, 11\}$ and the binary operation $+_{12}$. Among things we can do with **G** is that we can examine its subgroups.

```
Subgroups[G]
```

```
{Groupoid({0}, (#1 + #2) mod 12 &), Groupoid({0, 2, 4, 6, 8, 10}, (#1 + #2) mod 12 &),  
Groupoid({0, 3, 6, 9}, (#1 + #2) mod 12 &), Groupoid({0, 4, 8}, (#1 + #2) mod 12 &),  
Groupoid({0, 6}, (#1 + #2) mod 12 &), Groupoid({0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11}, (#1 + #2) mod 12 &)}
```

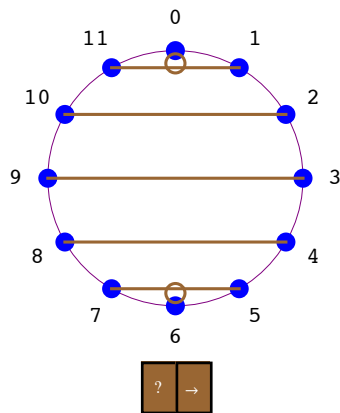
We can view the inverses of elements in a variety of ways. For example, we can get them paired up. Notice that two of the elements, 0 and 6 invert themselves.

```
Inverses[G]
```

$$\begin{pmatrix} 0 & 0 \\ 1 & 11 \\ 2 & 10 \\ 3 & 9 \\ 4 & 8 \\ 5 & 7 \\ 6 & 6 \end{pmatrix}$$

There is a "Visual Mode" that gives us a different view of the inverses. The boxes with "?" and "→" give further information in you are reading this in a *Mathematica* Notebook and have the package installed.

Inverses [**G**, **Mode** \rightarrow **Visual**]



The package was designed for teaching a first course in abstract algebra and so it has features that are more basic than other abstract algebra resources. For example, we can ask **G** is really a group and get quite a bit of information.

GroupQ[G, Mode → Textual]

Given a set S and an operation $*$, we call the pair $(S, *)$ a group if S is closed under the operation $*$, there is an identity element, every element has an inverse and the operation $*$ is associative.

We say a Groupoid G has an identity e if for all other elements g in G we have $e + g = g + e = g$ (where $+$ indicates the operation). In this case, $\mathbb{Z}[12]$ has the identity 0.



We say that a set S is closed under an operation op if whenever we have x and y in S , we also have $op[x,y]$ (or $x \sim op \sim y$) in S . In this case, the Groupoid $\mathbb{Z}[12]$ is indeed closed.



Given a Groupoid G , we say an element g in G has an inverse h if G has an identity, say e , and $g + h = h + g = e$ (where $+$ indicates the operation). The Groupoid $\mathbb{Z}[12]$ has an inverse for every element. Here they are:

x	x^{-1}
0	0
1	11
2	10
3	9
4	8
5	7
6	6



Given a structured set S (Groupoid or Ringoid), we say the operation $*$ is associative if for every g, h , and k in S we have $(g*h)*k = g*(h*k)$, where $*$ is the group operation. In this case, $\mathbb{Z}[12]$ is associative. Consider the following table illustrating random triples that associate. Pay attention to the last two columns.

i	j	k	$(i*j)*k$	$i*(j*k)$
2	11	4	5	5
3	1	4	8	8
8	9	8	1	1
7	5	2	2	2
8	8	2	6	6
10	6	9	1	1
4	9	7	8	8
4	5	2	11	11
11	4	1	4	4
5	10	6	9	9



This package also has much more capabilities than what we've covered so far and we will revisit it in Chapters 15 and 16.

Groups in Sage

Abstract Algebra seems to have been given a much higher priority in the design of Sage than it was in *Mathematica*. Again, the capabilities far exceed what we've touched on in the theory, but here are a few examples that you should understand. Here is how to generate the group \mathbb{Z}_{14} .

```
G=AbelianGroup(1,[14])
```

```
G.list()
[1, f, f^2, f^3, f^4, f^5, f^6, f^7, f^8, f^9, f^10, f^11, f^12, f^13]
```

There is no output from assigning G. The elements of G are generated from the `list` method. The connection with \mathbb{Z}_{14} is that when we multiply powers of `f`, the exponents are added with $+_{14}$. Among other things we can ask whether G is abelian and what its subgroups are.

```
G.is_abelian()
True
G.subgroups()
[Multiplicative Abelian Group isomorphic to C2 x C7, which is the
subgroup of
Multiplicative Abelian Group isomorphic to C14
generated by [f], Multiplicative Abelian Group isomorphic to C7, which
is the subgroup of
Multiplicative Abelian Group isomorphic to C14
generated by [f^2], Multiplicative Abelian Group isomorphic to C2, which
is the subgroup of
Multiplicative Abelian Group isomorphic to C14
generated by [f^7], Trivial Abelian Group, which is the subgroup of
Multiplicative Abelian Group isomorphic to C14
generated by []]
```

SUPPLEMENTARY EXERCISES FOR CHAPTER 11

Section 11.1

1. $V = \{a, b, c\}$ is a set with operations $+$ and \cdot defined by the following "addition" and "multiplication" tables:

$+$	a	b	c
a	a	b	c
b	b	c	a
c	c	a	b

\cdot	a	b	c
a	a	a	a
b	a	b	c
c	a	c	b

- With respect to V under $+$ determine,
 - The identity (i.e., the "zero" of the addition).
 - The inverse of each element, that is, $-a$, $-b$, and $-c$.
 - With respect to V under \cdot determine,
 - The identity (i.e., the "one" of the multiplication).
 - The inverse of each element different from "zero."
 - Is $+$ distributive over \cdot ? Is \cdot distributive over $+$?
2. (a) Determine whether the following are valid binary operations on the given sets. Explain fully.
- Matrix addition on $A = \left\{ \begin{pmatrix} a & b \\ c & 0 \end{pmatrix} \mid a, b, c \in \mathbb{R} \right\}$
 - Matrix multiplication on the set A above.
 - On \mathbb{Q}^+ , define $*$ by $a * b = (a \cdot b)/2$.
 - Function composition on $A^A = \{f: A \rightarrow A\}$, where A is $\{1, 2, 3\}$.
 - Function composition on $B = \{f \in A^A \mid f \text{ is a bijection}\}$.
- (b) For each binary operation above give the identity element if it exists. Explain.
- (c) Determine which of the above binary operations are commutative and which are associative.
3. Let S = set of all bijections of a set A , and let \circ be function composition. Does \circ have the inverse property? Does function composition have the involution property? Explain.
4. Does $+$ on $M_{2 \times 2}(\mathbb{R})$ have the inverse property? Does $+$ have the involution property? Explain.
5. Prove that the odd integers are closed under multiplication but not under addition. Are the even integers closed under both addition and multiplication? Prove your answers.

Section 11.2

6. (a) Show that \mathbb{R}^2 is a group under componentwise addition, that is,
- $$(a_1, a_2) + (b_1, b_2) = (a_1 + a_2, b_1 + b_2).$$
- (b) Show that $\{(x, 2x) \mid x \in \mathbb{R}\}$ is a group under componentwise addition. Draw the graph of this subset. Describe similar subsets of \mathbb{R}^2 that are also groups.
7. Prove that the set of all 2×2 invertible matrices (over \mathbb{R}) is a group under matrix multiplication. Assume, as indicated in Chapter 5, that the associative law is true for matrices under multiplication. This group is called the *general linear group* of degree 2 over \mathbb{R} , and it is denoted by $GL(2, \mathbb{R})$. It is given this name because these matrices are matrix representations of linear motions of \mathbb{R}^2 .
8. Prove that $\left\{ A = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \mid \det A = 1 \right\}$ is a group under matrix multiplication. Assume that the associative law is true under matrix multiplication. This group is called the *special linear group* of degree 2 over \mathbb{R} and it is denoted $SL(2, \mathbb{R})$.
9. Show that \mathbb{R} is a group under the operation $*$ defined by $a * b = a + b + 5$ for $a, b \in \mathbb{R}$.
10. (a) let $B_{3 \times 3}$ be the set of all 3×3 Boolean (adjacency) matrices discussed in Section 6.4. Is $B_{3 \times 3}$ a monoid under Boolean addition? Is it a group? Explain.
- (b) Is $B_{3 \times 3}$ a monoid under Boolean multiplication? Is it a group? Explain.

Section 11.3

11. Define $*$ on \mathbb{Q}^+ by $a * b = (a \cdot b)/2$. Prove that $[\mathbb{Q}^+; *]$ is a group.
12. Let G be the group \mathbb{R} under the operation $a * b = a + b + 5$ for $a, b \in \mathbb{R}$. Solve the following equations for x in G .

- (a) $x * 3 = 5$ (d) $x^2 = 2$
 (b) $2 * x * 4 = 6$ (e) $4 * x^2 = 5$
 (c) $x^3 = 7$

13. Solve the equation $A * X * B = C$ in $GL(2, \mathbb{R})$ where

$$A = \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}, \quad B = \begin{pmatrix} 4 & 2 \\ 2 & 2 \end{pmatrix}, \quad \text{and} \quad C = \begin{pmatrix} 2 & 1 \\ 0 & 1 \end{pmatrix}.$$

14. Prove that if $[G; *]$ is a group, $(a * b)^n = a^n * b^n$ for all $n \geq 1$ and $a, b \in G$ if and only if $[G; *]$ is an abelian group.

Section 11.4

15. Calculate the following in \mathbb{Z}_5 :

- (a) $3 +_5 8$
 (b) $(-3) \times_5 2$
 (c) $(3 \times_5 2) +_5 (2 \times_5 2)$
 (d) 2^{-1} (i.e., the multiplicative inverse of 2)

16. (a) Prove that $\{1, 3, 5, 7\}$, is a group under \times_8 . Write out its group table.

- (b) Let $U(\mathbb{Z}_n)$ stand for the elements of \mathbb{Z}_n , which have inverses under \times_n . Convince yourself that $U(\mathbb{Z}_n)$ is a group under \times_n .

- (c) Prove that the elements of $U(\mathbb{Z}_n)$ are those elements $a \in \mathbb{Z}_n$ such that $\gcd(a, n) = 1$. You may use the fact that $\gcd(a, b) = 1 \Leftrightarrow$ there exist integers s and t such that $sa + tb = 1$.

Section 11.5

17. (a) Recall from "Supplementary Exercises," Section 11.4, that $U(\mathbb{Z}_8)$ is a group under \times_8 . List all cyclic subgroups of this group.

- (b) Is $U(\mathbb{Z}_8)$ a cyclic group? Explain.

18. (a) Use Theorem 11.5.1 to prove that the set of even integers is a subgroup of the group \mathbb{Z} (under $+$).

- (b) Is the set of odd integers a subgroup of the group \mathbb{Z} (under $+$)?

19. Prove that $SL(2, \mathbb{R})$ is a subgroup of $GL(2, \mathbb{R})$. See Exercises 7 and 8 above for an explanation of this notation.

20. Recall that $M_{2 \times 2}(\mathbb{R})$ is a group under addition.

- (a) Is $A = \left\{ \begin{pmatrix} a & b \\ a & 0 \end{pmatrix} \mid a, b \in \mathbb{R} \right\}$ a subgroup of $M_{2 \times 2}(\mathbb{R})$?

- (b) Is $B = \left\{ \begin{pmatrix} a & b \\ b & 1 \end{pmatrix} \mid a, b \in \mathbb{R} \right\}$ a subgroup of $M_{2 \times 2}(\mathbb{R})$?

- (c) Are either of the subsets in parts a and b subgroups of $GL(2, \mathbb{R})$?

21. Let $B_{3 \times 3}$ be the monoid of all 3×3 Boolean matrices, under Boolean addition. Let S be a subset of $B_{3 \times 3}$ consisting of all 3×3 matrices that represent symmetric relations. Is S a submonoid of $B_{3 \times 3}$?

Section 11.6

22. Using the data structure in the text for doubly linked lists with six-bit addresses, what are the addresses of the records containing A and D? Write your answer as a sum in the group \mathbb{Z}_2^6 and then as an address.

?	011100	000011	?
A	B	C	D
	010101	001011	

23. Determine the inverse of each element in the respective group.

- (a) $(2, 3, 5)$ in $\mathbb{Z}_3 \times \mathbb{Z}_7 \times \mathbb{Z}_{25}$
 (b) $(1, 0, 1, 1)$ in \mathbb{Z}^4

(c) $(3, 2)$ in $\mathbb{R}^+ \times \mathbb{Z}_6$

(d) $(2, 3, 5)$ in \mathbb{R}^3

24. Determine the identity elements in the following groups:

(a) $\mathbb{R}^+ \times \mathbb{R}^+$

(b) $\mathbb{R}^+ \times \mathbb{Z}_3$

(c) $\text{GL}(2, \mathbb{R}) \times \mathbb{R}^3$

25. Which of the following groups are abelian? Explain.

(a) $\mathbb{Z}_2 \times \mathbb{Z}_{24} \times \mathbb{Z}_{75}$

(b) $\text{GL}(2, \mathbb{R}) \times \mathbb{Z}_2$

(c) \mathbb{Z}^n

26. Is $\{0, 3\} \times \{0, 4, 8\}$ a subgroup of $\mathbb{Z}_6 \times \mathbb{Z}_8$? Explain.

Section 11.7

27. Prove that the cyclic subgroup $\langle 4 \rangle$ of \mathbb{Z}_{16} is isomorphic to \mathbb{Z}_4 .

28. Let $G = \{ \&, \$, \% \}$. Given that $[G; *]$ is a group and that it is isomorphic to the group $[\mathbb{Z}_3; +_3]$ with isomorphism $T : G \rightarrow \mathbb{Z}_3$ defined by $T(\&) = 1$, $T(\$) = 2$, and $T(\%) = 0$. What are

(a) $\$ * \$$ (b) The identity of $[G; *]$

29. Let U be a set and $\mathcal{P}_U = \{\text{propositions over the set } U\}$. It can be shown that the algebraic system $[\mathcal{P}_U; \sim, \wedge, \vee]$ is isomorphic to $[\mathcal{P}(U); \cdot, \cap, \cup]$.

(a) Explain what this means.

(b) How does this help you understand the language of the algebra of propositions?

(c) Give the "propositional" analogue to the following statement: If $A \cap B^c = \emptyset$ and $A \cap B = \emptyset$ then $A = \emptyset$.

30. Write out the operation tables for the following systems:

(a) $\{0, 1\}; +, \cdot$ where $+$ and \cdot denote Boolean addition and multiplication.

(b) $\{-1, 1\}; \wedge, \vee$ where $i \wedge j$ and $i \vee j$ denote the largest and smallest, respectively, of i and j .

(c) $[\mathbb{Z}_2; +_2, \times_2]$.

Are these systems isomorphic? Explain.

31. Prove that the group \mathbb{C} , under $+$, is isomorphic to the group \mathbb{R}^2 , under $+$.

32. Determine which of the following groups are isomorphic. Explain.

(a) \mathbb{R}_3 , the 3×3 rook matrices, and \mathbb{Z}_6

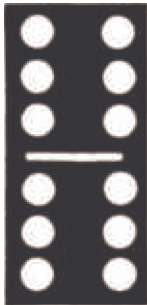
(b) \mathbb{R}_3 and $S_A = \{f \in A^A : f \text{ is a bijection}\}$, where A is $\{1, 2, 3\}$.

(c) \mathbb{Z}_6 and $U(\mathbb{Z}_7)$

33. Prove that \mathbb{R}^4 under addition, is isomorphic to $M_{2 \times 2}(\mathbb{R})$, under addition.

34. Prove that the group $[U(\mathbb{Z}_8); \times_8]$ is isomorphic to $[\mathbb{Z}_4; +_4]$.

chapter 12



MORE MATRIX ALGEBRA

GOALS

In Chapter 5 we studied matrix operations and the algebra of sets and logic. We also made note of the strong resemblance of matrix algebra to elementary algebra. The reader should briefly review this material. In this chapter we shall look at a powerful matrix tool in the applied sciences—namely, a technique for solving systems of linear equations. We will then use this process for determining the inverse of $n \times n$ matrices, $n \geq 2$, when they exist. We conclude by a development of the diagonalization process, with a discussion of several of its applications.

12.1 Systems of Linear Equations

The method of solving systems of equations by matrices that we will look at is based on procedures involving equations that we are familiar with from previous mathematics courses. The main idea is to reduce a given system of equations to another simpler system that has the same solutions.

Definition: Solution Set. Given a system of equations involving real variables x_1, x_2, \dots, x_n , the solution set of the system is the set of n -tuples in \mathbb{R}^n , (a_1, a_2, \dots, a_n) such that the substitutions $x_1 = a_1, x_2 = a_2, \dots, x_n = a_n$ make all the equations true.

In general, if the variables are from a set S , then the solution set will be a subset of S^n . For example, in number theory mathematicians study Diophantine equations, where the variables can only take on integer values instead of real values.

Definition: Equivalent Systems of Equations. Two systems of linear equations are called equivalent if they have the same set of solutions.

Example 12.1.1. The previous definition tells us that if we know that the system

$$\begin{aligned} 4x_1 + 2x_2 + x_3 &= 1 \\ 2x_1 + x_2 + x_3 &= 4 \\ 2x_1 + 2x_2 + x_3 &= 3 \end{aligned}$$

is equivalent to the system

$$\begin{aligned} x_1 + 0x_2 + 0x_3 &= -1 \\ 0x_1 + x_2 + 0x_3 &= -1 \\ 0x_1 + 0x_2 + x_3 &= 7 \end{aligned}$$

then both systems have the solution set $\{(-1, -1, 7)\}$. In other words, the values $x_1 = -1$, $x_2 = -1$, and $x_3 = 7$ are the only values of the variables that make all three equations in either system true.

Theorem 12.1.1. Elementary Operations on Equations. If any sequence of the following operations is performed on a system of equations, the resulting system is equivalent to the original system:

- (1) Interchange any two equations in the system.
- (2) Multiply both sides of any equation by a nonzero constant.

(3) Multiply both sides of any equation by a nonzero constant and add the result to a second equation in the system, with the sum replacing the latter equation.

Let us now use the above theorem to work out the details of Example 12.1.1 and see how we can arrive at the simpler system..

Step 1. We will first change the coefficient of x_1 in the first equation to one and then use it as a pivot to obtain 0's for the coefficients of x_1 in Equations 2 and 3.

$$\begin{array}{lcl} & 4x_1 + 2x_2 + x_3 = 1 & \\ (1.1) & 2x_1 + x_2 + x_3 = 4 & \text{Multiply Equation 1 by } \frac{1}{4} \text{ to obtain} \\ & 2x_1 + 2x_2 + x_3 = 3 & \end{array}$$

$$\begin{array}{lcl} & x_1 + \frac{x_2}{2} + \frac{x_3}{4} = \frac{1}{4} & \\ (1.2) & 2x_1 + x_2 + x_3 = 4 & \text{Multiply Equation 1 by } -2 \text{ and} \\ & 2x_1 + 2x_2 + x_3 = 3 & \end{array}$$

add the result to Equation 3 to obtain

$$\begin{array}{lcl} & x_1 + \frac{x_2}{2} + \frac{x_3}{4} = \frac{1}{4} & \\ (1.3) & 0x_1 + 0x_2 + \frac{x_3}{2} = \frac{7}{2} & \text{Multiply Equation 1 by } -2 \text{ and add} \\ & 2x_1 + 2x_2 + x_3 = 3 & \end{array}$$

the result to Equation 3 to obtain

$$\begin{array}{lcl} & x_1 + \frac{x_2}{2} + \frac{x_3}{4} = \frac{1}{4} & \\ (1.4) & 0x_1 + 0x_2 + \frac{x_3}{2} = \frac{7}{2} & \\ & 0x_1 + x_2 + \frac{x_3}{2} = \frac{5}{2} & \end{array}$$

Note: We've explicitly written terms with zero coefficients such as $0x_1$ to make a point that all variables can be thought of as being involved in all equations. After this example we will discontinue this practice in favor of the normal practice of making these terms "disappear."

Step 2. We would now like to proceed in a fashion analogous to Step 1—namely, multiply the coefficient of x_2 in the second equation by a suitable number so that the result is 1. Then use it as a pivot to obtain 0's as coefficients for x_2 in the first and third equations. This is clearly impossible (Why?), so we will first interchange Equations 2 and 3 and proceed as outlined above.

$$\begin{array}{lcl} & x_1 + \frac{x_2}{2} + \frac{x_3}{4} = \frac{1}{4} & \\ (2.1) & 0x_1 + 0x_2 + \frac{x_3}{2} = \frac{7}{2} & \text{Interchange Equations 2 and 3 to obtain} \\ & 0x_1 + x_2 + \frac{x_3}{2} = \frac{5}{2} & \end{array}$$

$$\begin{array}{lcl} & x_1 + \frac{x_2}{2} + \frac{x_3}{4} = \frac{1}{4} & \\ (2.2) & 0x_1 + x_2 + \frac{x_3}{2} = \frac{5}{2} & \text{Multiply Equation 2 by } -\frac{1}{2} \text{ and add} \\ & 0x_1 + 0x_2 + \frac{x_3}{2} = \frac{7}{2} & \text{the result to Equation 1 to obtain} \end{array}$$

$$\begin{array}{lcl} & x_1 + 0x_2 + 0x_3 = -1 & \\ (2.3) & 0x_1 + x_2 + \frac{x_3}{2} = \frac{5}{2} & \\ & 0x_1 + 0x_2 + \frac{x_3}{2} = \frac{7}{2} & \end{array}$$

Step 3. Next, we will change the coefficient of x_3 in the third equation to one and then use it as a pivot to obtain 0's for the coefficients of x_3 in Equations 1 and 2.

$$\begin{array}{lcl} & x_1 + 0x_2 + 0x_3 = -1 & \\ (3.1) & 0x_1 + x_2 + \frac{x_3}{2} = \frac{5}{2} & \text{Multiply Equation 3 by 2 to obtain} \\ & 0x_1 + 0x_2 + x_3 = 7 & \end{array}$$

$$\begin{array}{lcl} & x_1 + 0x_2 + 0x_3 = -1 & \\ (3.2) & 0x_1 + x_2 + \frac{x_3}{2} = \frac{5}{2} & \text{Multiply Equation 3 by } -\frac{1}{2} \text{ and add the result} \\ & 0x_1 + 0x_2 + x_3 = 7 & \text{to Equation 2 to obtain} \end{array}$$

$$\begin{array}{lcl} & x_1 + 0x_2 + 0x_3 = -1 & \\ (3.3) & 0x_1 + x_2 + 0x_3 = -1 & \\ & 0x_1 + 0x_2 + x_3 = 7 & \end{array}$$

From the system of equations in Step 3.3, we see that the solution to the original system (Step 1.1) is $x_1 = -1$, $x_2 = -1$, and $x_3 = 7$.

In the above sequence of steps, we note that the variables serve the sole purpose of keeping the coefficients in the appropriate location. This we can effect by using matrices. The matrix of the system given in Step 1.1 is

$$\begin{pmatrix} 4 & 2 & 1 & 1 \\ 2 & 1 & 1 & 4 \\ 2 & 2 & 1 & 3 \end{pmatrix}$$

where the matrix of the first three columns is called the coefficient matrix and the complete matrix is referred to as the augmented matrix. Since we are now using matrices to solve the system, we will translate Theorem 12.1.1 into matrix language.

Definition: Elementary Row Operations. The following operations on a matrix are called elementary row operations:

- (1) Interchange any two rows of the matrix.
- (2) Multiply any row of the matrix by a nonzero constant.
- (3) Multiply any row of the matrix by a nonzero constant and add the result to a second row, with the sum replacing the second row.

Definition: Row Equivalent. Two matrices, A and B , are said to be row-equivalent if one can be obtained from the other by any one elementary row operation or by any sequence of elementary row operations.

If we use the notation R_i to stand for Row i of a matrix and \longrightarrow to stand for row equivalence, then

$$A \xrightarrow{cR_i + R_j} B$$

means that the matrix B is obtained from the matrix A by multiplying the Row i of A by c and adding the result to Row j . The operation of multiplying row i by c is indicated by

$$A \xrightarrow{cR_i} B$$

while interchanging rows i and j is denoted by

$$A \xrightarrow{R_i \leftrightarrow R_j} B.$$

The matrix notation for the system given in Step 1.1 with the subsequent steps are:

$$\begin{aligned} \begin{pmatrix} 4 & 2 & 1 & 1 \\ 2 & 1 & 1 & 4 \\ 2 & 2 & 1 & 3 \end{pmatrix} &\xrightarrow{\frac{1}{4}R_1} \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ 2 & 1 & 1 & 4 \\ 2 & 2 & 1 & 3 \end{pmatrix} \xrightarrow{-2R_1 + R_2} \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ 0 & 0 & \frac{1}{2} & \frac{7}{2} \\ 2 & 2 & 1 & 3 \end{pmatrix} \\ &\xrightarrow{-2R_1 + R_3} \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ 0 & 0 & \frac{1}{2} & \frac{7}{2} \\ 0 & 1 & \frac{1}{2} & \frac{5}{2} \end{pmatrix} \xrightarrow{R_2 \leftrightarrow R_3} \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{4} & \frac{1}{4} \\ 0 & 1 & \frac{1}{2} & \frac{5}{2} \\ 0 & 0 & \frac{1}{2} & \frac{7}{2} \end{pmatrix} \\ &\xrightarrow{-\frac{1}{2}R_2 + R_1} \begin{pmatrix} 1 & 0 & 0 & -1 \\ 0 & 1 & \frac{1}{2} & \frac{5}{2} \\ 0 & 0 & \frac{1}{2} & \frac{7}{2} \end{pmatrix} \xrightarrow{2R_3} \begin{pmatrix} 1 & 0 & 0 & -1 \\ 0 & 1 & \frac{1}{2} & \frac{5}{2} \\ 0 & 0 & 1 & 7 \end{pmatrix} \\ &\xrightarrow{-\frac{1}{2}R_3 + R_2} \begin{pmatrix} 1 & 0 & 0 & -1 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & 7 \end{pmatrix} \end{aligned}$$

This again gives us the solution. This procedure is called the *Gauss-Jordan elimination method*.

It is important to remember when solving any system of equations via this or any similar approach that at any step in the procedure we can rewrite the matrix in "equation format" to help us to interpret the meaning of the augmented matrix.

In Example 12.1.1 we obtained a unique solution, only one triple, namely $(-1, -1, 7)$, which satisfies all three equations. For a system involving three unknowns, are there any other possible results? To answer this question, let's review some basic facts from analytic geometry.

The graph of a linear equation in three-dimensional space is a plane. So geometrically we can visualize the three linear equations as three planes in three-space. Certainly the three planes can intersect in a unique point, as in Example 12.1.1, or two of the planes could be parallel. If two planes are parallel, there are no common points of intersection; that is, there are no triple of real numbers that will satisfy all three equations. Also, the three planes could intersect along a common axis or line. In this case, there would be an infinite number of real number triples in \mathbb{R}^3 that would satisfy all three equations. Finally if all three equations describe the same plane, the solution set would be that plane. We generalize;

In a system of n linear equations, n unknowns, there can be:

- (1) a unique solution,
- (2) no solution, or
- (3) an infinite number of solutions.

To illustrate these points, consider the following examples:

Example 12.1.2. Find all solutions to the system

$$x_1 + 3x_2 + x_3 = 2$$

$$x_1 + x_2 + 5x_3 = 4$$

$$2x_1 + 2x_2 + 10x_3 = 6$$

The reader can verify that the augmented matrix of this system,

$$\left(\begin{array}{ccc|c} 1 & 3 & 1 & 2 \\ 1 & 1 & 5 & 4 \\ 2 & 2 & 10 & 6 \end{array} \right),$$

reduces to

$$\left(\begin{array}{ccc|c} 1 & 3 & 1 & 2 \\ 1 & 1 & 5 & 4 \\ 0 & 0 & 0 & -2 \end{array} \right) \quad (\text{See exercise 4 of this section.})$$

We can row-reduce this matrix further if we wish. However, any further row-reduction will not substantially change the last row, which, in equation form, is $0x_1 + 0x_2 + 0x_3 = -2$, or simply $0 = -2$. It is clear that we cannot find real numbers x_1 , x_2 , and x_3 that will satisfy this equation, hence we cannot find real numbers that will satisfy all three original equations simultaneously. When this occurs, we say that the system has no solution, or the solution set is empty.

Example 12.1.3. Next let's attempt to find all of the solutions to:

$$x_1 + 6x_2 + 2x_3 = 1$$

$$2x_1 + x_2 + 3x_3 = 2$$

$$4x_1 + 2x_2 + 6x_3 = 4$$

The augmented matrix for the system,

$$\left(\begin{array}{ccc|c} 1 & 6 & 2 & 1 \\ 2 & 1 & 3 & 2 \\ 4 & 2 & 6 & 4 \end{array} \right)$$

reduces to

$$\left(\begin{array}{ccc|c} 1 & 0 & \frac{16}{11} & 1 \\ 0 & 1 & \frac{1}{11} & 0 \\ 0 & 0 & 0 & 0 \end{array} \right)$$

If we apply additional elementary row operations to this matrix, it will only become more complicated. In particular, we cannot get a one in the third row, third column. Since the matrix is in simplest form, we will express it in equation format to help us determine the solution set.

$$x_1 + \frac{16}{11}x_3 = 1$$

$$x_2 + \frac{1}{11}x_3 = 0$$

$$0 = 0$$

Any real numbers will satisfy the last equation. However, the first equation can be rewritten as $x_1 = 1 - \frac{16}{11}x_3$, which describes the coordinate x_1 in terms of x_3 . Similarly, the second equation gives x_2 in terms of x_3 . A convenient way of listing the solutions of this system is to use set notation. If we call the solution set of the system S , then

$$S = \left\{ \left(1 - \frac{16}{11}x_3, -\frac{1}{11}x_3, x_3 \right) \mid x_3 \in \mathbb{R} \right\}.$$

What this means is that if we wanted to list all solutions, we would replace x_3 by all possible numbers. Clearly, there is an infinite number of solutions, two of which are $(1, 0, 0)$ and $(-15, -1, 11)$.

A Word Of Caution: Frequently we may obtain “different-looking” answers to the same problem when a system has an infinite number of answers. Assume a student's solutions set to Example 12.1.3 is $A = \{(1 + 16x_2, x_2, -11x_3) \mid x_3 \in \mathbb{R}\}$. Certainly the result described by S looks different from that described by A . To see whether they indeed describe the same set, we wish to determine whether every solution produced in S can be generated in A . For example, the solution generated by S when $x_3 = 11$ is $(-15, -1, 11)$. The same triple can be produced by A by taking $x_2 = -1$. We must prove that every solution described in S is described in A and, conversely, that every solution described in A is described in S . (See Exercise 6 of this section.)

To summarize the procedure in the Gauss-Jordan technique for solving systems of equations, we attempt to obtain 1's along the main diagonal of the coefficient matrix with 0's above and below the diagonal, as in Example 12.1.1. We may find in attempting this that the closest we can come is to put the coefficient matrix in "simplest" form, as in Example 12.1.3, or we may find that the situation of Example 12.1.1 evolves as part of the process. In this latter case, we can terminate the process and state that the system has no solutions. The final matrix forms of Examples 12.1.1 and 12.1.3 are called echelon forms.

In practice, larger systems of linear equations are solved using computers. Generally, the Gauss-Jordan algorithm is the most useful; however, slight variations of this algorithm are also used. The different approaches share many of the same advantages and disadvantages. The two major concerns of all methods are:

- (1) minimizing inaccuracies due to rounding off errors, and
- (2) minimizing computer time.

The accuracy of the Gauss-Jordan method can be improved by always choosing the element with the largest absolute value as the pivot element, as in the following algorithm.

Algorithm 12.1.1. Given a matrix equation $Ax = b$, where A is $n \times m$, let C be the augmented matrix $[A \mid b]$. The process of **row-reducing to echelon form** involves performing the following algorithm where $C_i =$ the i^{th} row of C :

```

i = 1
j = 1
while (i ≤ n and j ≤ m):
    # Find pivot in column j, starting in row i:
    maxi = i
    for k = i+1 to n:
        if abs(C[k,j]) > abs(C[maxi,j]) then
            maxi := k
    if C[maxi,j] ≠ 0 then
        interchange rows i and maxi
        divide each entry in row i by C[i,j]
        # Now C[i,j] will have the value 1.
        for u = i+1 to n:
            subtract C[u,j] * Ci from Cu
            # Now C[u,j] will be 0
        i := i + 1
    end if
    j = j + 1
end while

```

At the end of this algorithm, with the final form of C you can revert back to the equation form of the system and a solution should be clear. In general,

(a) If any row of C is all zeros, it can be ignored.

(b) If any row of C has all zero entries except for the entry in the $(m+1)^{\text{st}}$ position, the system has no solution. Otherwise, if a column has no pivot, the variable corresponding to it is a **free variable**. Variables corresponding to pivots are **basic variables** and can be expressed in terms of the free variables.

Example 12.1.4. If we apply Algorithm 12.1.1 to the system

$$\begin{aligned} 5x_1 + x_2 + 2x_3 + x_4 &= 2 \\ 3x_1 + x_2 - 2x_3 &= 5 \\ x_1 + x_2 + 3x_3 - x_4 &= -1 \end{aligned}$$

the augmented matrix

$$C = \begin{pmatrix} 5 & 1 & 2 & 1 & 2 \\ 3 & 1 & -2 & 0 & 5 \\ 1 & 1 & 3 & -1 & -1 \end{pmatrix}$$

is reduced to a new value of C :

$$C = \begin{pmatrix} 1 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 1 & 0 & -\frac{3}{2} & \frac{3}{2} \\ 0 & 0 & 1 & 0 & -1 \end{pmatrix}$$

therefore x_4 is a free variable in the solution and general solution of the system is

$$x = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} = \begin{pmatrix} \frac{1}{2} - \frac{1}{2}x_4 \\ \frac{3}{2} + \frac{3}{2}x_4 \\ -1 \\ x_4 \end{pmatrix}$$

This conclusion is easy to see if you revert back to the equations that the final value matrix C represents.



Mathematica Note

The *Mathematica* function **RowReduce** does the same reduction as described in Algorithm 12.1.1. For example, here is the result for the system in Example 12.1.4.

$$\text{RowReduce}\left[\begin{pmatrix} 5 & 1 & 2 & 1 & 2 \\ 3 & 1 & -2 & 0 & 5 \\ 1 & 1 & 3 & -1 & -1 \end{pmatrix}\right]$$

$$\begin{pmatrix} 1 & 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 1 & 0 & -\frac{3}{2} & \frac{3}{2} \\ 0 & 0 & 1 & 0 & -1 \end{pmatrix}$$

Options[RowReduce]

{Method → Automatic, Modulus → 0, Tolerance → Automatic, ZeroTest → Automatic}

Only one caution: One needs to be aware that if the pivoting process continues into the last column, which *Mathematica* will do, there will not be a solution to the system. For example the system

$$\begin{aligned} 2x_1 - x_2 &= 1 \\ 3x_2 - x_1 &= 5 \\ x_1 + 5x_2 &= 7 \end{aligned}$$

has augmented matrix

$$C = \begin{pmatrix} 2 & -1 & 1 \\ -1 & 3 & 5 \\ 1 & 5 & 7 \end{pmatrix}.$$

Here is the computation to row-reduce:

$$\text{RowReduce}\left[\begin{pmatrix} 2 & -1 & 1 \\ -1 & 3 & 5 \\ 1 & 5 & 7 \end{pmatrix}\right]$$

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

The last row of the final form of C is $0 = 1$ and so there is no solution to the original system.



Sage Note

Given an augmented matrix, C , there is a matrix method called `echelon_form` that can be used to row reduce C . Here is the result for the system in Example 12.1.4. In the assignment of a matrix value to C , notice that the first argument is `QQ`, which indicates that the entries should be rational numbers. As long as all the entries are rational, which is the case here since integers are rational, the row-reduced matrix will be all rational.

```
C = Matrix(QQ, [[5,1,2,1,2],[3,1,-2,0,5],[1,1,3,-1,-1]])
C.echelon_form()
[ 1  0  0  1/2  1/2 ]
[ 0  1  0 -3/2  3/2 ]
[ 0  0  0  1   0  -1 ]
```

If we didn't specify the set from which entries are taken, it would assumed to be the integers and we would not get a fully row-reduced matrix. The next step would involve multiplying row 3 by $\frac{1}{9}$, which isn't an integer.

```
C2 = Matrix([[5,1,2,1,2],[3,1,-2,0,5],[1,1,3,-1,-1]])
C2.echelon_form()
```

$$\begin{bmatrix} 1 & 1 & 3 & -1 & -1 \\ 0 & 2 & 2 & -3 & 1 \\ 0 & 0 & 9 & 0 & -9 \end{bmatrix}$$

This is why we would avoid specifying real entries:

```
C3 = Matrix(RR, [[5,1,2,1,2],[3,1,-2,0,5],[1,1,3,-1,-1]])
C3.echelon_form()
[ 1.000000000000000 0.000000000000000 0.000000000000000 0.500000000000000 0.500000000000000]
[ 0.000000000000000 1.000000000000000 0.000000000000000 -1.500000000000000 1.500000000000000]
[ 0.000000000000000 0.000000000000000 1.000000000000000 4.93432455388958e-17 -1.000000000000000]
```

This is the default number of decimal places, which could be controled and the single small number in row three column four isn't exactly zero because of round-off and we could just set it to zero. However, the result isn't as nice and clean as the rational output in this case.

EXERCISES FOR SECTION 12.1

A Exercises

1. Solve the following systems by describing the solution sets completely:

- (a) $2x_1 + x_2 = 3$
 $x_1 - x_2 = 1$
 $2x_1 + x_2 + 3x_3 = 5$
- (b) $4x_1 + x_2 + 2x_3 = -1$
 $8x_1 + 2x_2 + 4x_3 = -2$
 $x_1 + x_2 + 2x_3 = 1$
- (c) $x_1 + 2x_2 - x_3 = -1$
 $x_1 + 3x_2 + x_3 = 5$
- (d) $x_1 - x_2 + 3x_3 = 7$
 $x_1 + 3x_2 + x_3 = 4$

2. Solve the following systems by describing the solution sets completely:

- (a) $2x_1 + 2x_2 + 4x_3 = 2$
 $2x_1 + x_2 + 4x_3 = 0$
 $3x_1 + 5x_2 + x_3 = 0$
 $2x_1 + x_2 + 3x_3 = 2$
- (b) $4x_1 + x_2 + 2x_3 = -1$
 $8x_1 + 2x_2 + 4x_3 = 4$
 $x_1 + x_2 + 2x_3 + x_4 = 3$
- (c) $x_1 - x_2 + 3x_3 - x_4 = -2$
 $3x_1 + 3x_2 + 6x_3 + 3x_4 = 9$
 $6x_1 + 7x_2 + 2x_3 = 3$
- (d) $4x_1 + 2x_2 + x_3 = -2$
 $6x_1 + x_2 + x_3 = 1$
 $x_1 + x_2 - x_3 + 2x_4 = 1$
- (e) $x_1 + 2x_2 + 3x_3 + x_4 = 5$
 $x_1 + 3x_2 + 2x_3 - x_4 = -1$

3. Given that the final augmented matrices below obtained from Algorithm 12.1.1, identify the solutions sets. Identify the basic and free variables, and describe the solution set of the original system.

(a) $\begin{pmatrix} 1 & 0 & -5 & 0 & 1.2 \\ 0 & 1 & 4 & 0 & 2.6 \\ 0 & 0 & 0 & 1 & 4.5 \end{pmatrix}$ (c) $\begin{pmatrix} 1 & 0 & 9 & 3 \\ 0 & 1 & 0 & 4 \\ 0 & 0 & 0 & 1 \end{pmatrix}$

(b) $\begin{pmatrix} 1 & 0 & 6 & 5 \\ 0 & 1 & -2 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix}$ (d) $\begin{pmatrix} 1 & 0 & 0 & -3 & 1 \\ 0 & 1 & 0 & 2 & 2 \\ 0 & 0 & 1 & -1 & 1 \end{pmatrix}$

4. (a) Write out the details of Example 12.1.2.
 (b) Write out the details of Example 12.1.3.
 (c) Write out the details of Example 12.1.4.
5. Solve the following systems using only mod 5 arithmetic. Your solutions should be n – tuples from \mathbb{Z}_5 .
- (a) $2x_1 + x_2 = 3$
 $x_1 + 4x_2 = 1$ (compare your solution to the system in 5(a))
 $x_1 + x_2 + 2x_3 = 1$
- (b) $x_1 + 2x_2 + 4x_3 = 4$
 $x_1 + 3x_2 + 3x_3 = 0$
6. (a) Use the solution set S of Example 12.1.3 to list three different solutions to the given system. Then show that each of these solutions can be described by the set A of Example 12.1.3.
 (b) Prove that $S = A$.

B Exercise

7. Given a system of n linear equations in n unknowns in matrix form $Ax = b$, prove that if b is a matrix of all zeros, then the solution set of $Ax = b$ is a subgroup of \mathbb{R}^n .

12.2 Matrix Inversion

In Chapter 5 we defined the inverse of an $n \times n$ matrix. We noted that not all matrices have inverses, but when the inverse of a matrix exists, it is unique. This enables us to define the inverse of an $n \times n$ matrix A as the unique matrix B such that $AB = BA = I$, where I is the $n \times n$ identity matrix. In order to get some practical experience, we developed a formula that allowed us to determine the inverse of invertible 2×2 matrices. We will now use the Gauss-Jordan procedure for solving systems of linear equations to compute the inverses, when they exist, of $n \times n$ matrices, $n \geq 2$. The following procedure for a 3×3 matrix can be generalized for $n \times n$ matrices, $n \geq 2$.

Example 12.2.1. Given the matrix

$$A = \begin{pmatrix} 1 & 1 & 2 \\ 2 & 1 & 4 \\ 3 & 5 & 1 \end{pmatrix}$$

we want to find the matrix

$$B = \begin{pmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{pmatrix},$$

if it exists, such that (a) $AB = I$ and (b) $BA = I$. We will concentrate on finding a matrix that satisfies Equation (a) and then verify that B also satisfies Equation (b).

$$\begin{pmatrix} 1 & 1 & 2 \\ 2 & 1 & 4 \\ 3 & 5 & 1 \end{pmatrix} \begin{pmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

is equivalent to

$$\begin{pmatrix} x_{11} + x_{21} + 2x_{31} & x_{12} + x_{22} + 2x_{32} & x_{13} + x_{23} + 2x_{33} \\ 2x_{11} + x_{21} + 4x_{31} & 2x_{12} + x_{22} + 4x_{32} & 2x_{13} + x_{23} + 4x_{33} \\ 3x_{11} + 5x_{21} + x_{31} & 3x_{12} + 5x_{22} + x_{32} & 3x_{13} + 5x_{23} + x_{33} \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad (12.2.a)$$

By definition of equality of matrices, this gives us three systems of equations to solve. The augmented matrix of one of the 12.2a systems, the one equating the first columns of the two matrices is:

$$\begin{pmatrix} 1 & 1 & 2 & 1 \\ 2 & 1 & 4 & 0 \\ 3 & 5 & 1 & 0 \end{pmatrix} \quad (12.2.b)$$

Using the Gauss-Jordan technique of Section 12.1, we have:

$$\begin{aligned}
\begin{pmatrix} 1 & 1 & 2 & 1 \\ 2 & 1 & 4 & 0 \\ 3 & 5 & 1 & 0 \end{pmatrix} &\xrightarrow{-2R_1+R_2} \begin{pmatrix} 1 & 1 & 2 & 1 \\ 0 & -1 & 0 & -2 \\ 3 & 5 & 1 & 0 \end{pmatrix} \xrightarrow{-3R_1+R_3} \begin{pmatrix} 1 & 1 & 2 & 1 \\ 0 & -1 & 0 & -2 \\ 0 & 2 & -5 & -3 \end{pmatrix} \\
&\xrightarrow{-1R_2} \begin{pmatrix} 1 & 1 & 2 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & 2 & -5 & -3 \end{pmatrix} \xrightarrow{\substack{-R_2+R_1 \\ \text{and } -2R_2+R_3}} \begin{pmatrix} 1 & 0 & 2 & -1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & -5 & -7 \end{pmatrix} \\
&\xrightarrow{-\frac{1}{5}R_3} \begin{pmatrix} 1 & 0 & 2 & -1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 7/5 \end{pmatrix} \xrightarrow{-2R_3+R_1} \begin{pmatrix} 1 & 0 & 0 & -\frac{19}{5} \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & \frac{7}{5} \end{pmatrix}
\end{aligned}$$

So $x_{11} = -19/5$, $x_{21} = 2$ and $x_{31} = 7/5$, which gives us the first column of the matrix B . The matrix form of the system to obtain x_{12} , x_{22} , and x_{32} , the second column of B , is:

$$\begin{pmatrix} 1 & 1 & 2 & 0 \\ 2 & 1 & 4 & 1 \\ 3 & 5 & 1 & 0 \end{pmatrix} \quad (12.2.c)$$

which reduces to

$$\begin{pmatrix} 1 & 0 & 0 & \frac{9}{5} \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -\frac{2}{5} \end{pmatrix} \quad (12.2.d)$$

The critical idea to note here is that the coefficient matrix in 12.2c is the same as the matrix in 12.2b, hence the sequence of row operations that we used to reduce the matrix in 12.2b can be used to reduce the matrix in 12.2c. To determine the third column of B , we reduce

$$\begin{pmatrix} 1 & 1 & 2 & 0 \\ 2 & 1 & 4 & 0 \\ 3 & 5 & 1 & 1 \end{pmatrix}$$

to obtain $x_{13} = 2/5$, $x_{23} = 0$ and $x_{33} = -1/5$. Here again it is important to note that the sequence of row operations used to "solve" this system is exactly the same as those we used in the first system. Why not save ourselves a considerable amount of time and effort and solve all three systems simultaneously? This we can effect by augmenting the coefficient matrix by the identity matrix I . We then have

$$\begin{pmatrix} 1 & 1 & 2 & 1 & 0 & 0 \\ 2 & 1 & 4 & 0 & 1 & 0 \\ 3 & 5 & 1 & 0 & 0 & 1 \end{pmatrix} \xrightarrow{\substack{\text{Same sequence of row} \\ \text{operations as above}}} \begin{pmatrix} 1 & 0 & 0 & -\frac{19}{5} & \frac{9}{5} & \frac{2}{5} \\ 0 & 1 & 0 & 2 & -1 & 0 \\ 0 & 0 & 1 & \frac{7}{5} & -\frac{2}{5} & -\frac{1}{5} \end{pmatrix}$$

So that

$$B = \begin{pmatrix} -\frac{19}{5} & \frac{9}{5} & \frac{2}{5} \\ 2 & -1 & 0 \\ \frac{7}{5} & -\frac{2}{5} & -\frac{1}{5} \end{pmatrix}$$

The reader should verify that $BA = I$ so that $A^{-1} = B$.

As the following theorem indicates, the verification that $BA = I$ is not necessary. The proof of the theorem is beyond the scope of this text. The interested reader can find it in most linear algebra texts.

Theorem 12.2.1. Let A be an $n \times n$ matrix. If a matrix B can be found such that $AB = I$, then $BA = I$, so that $B = A^{-1}$. In fact, to find A^{-1} , we need only find a matrix B that satisfies one of the two conditions $AB = I$ or $BA = I$.

It is clear from Chapter 5 and our discussions in this chapter that not all $n \times n$ matrices have inverses. How do we determine whether a matrix has an inverse using this method? The answer is quite simple: the technique we developed to compute inverses is a matrix approach to solving several systems of equations simultaneously.

Example 12.2.2. The reader can verify that if

$$A = \begin{pmatrix} 1 & 2 & 1 \\ -1 & -2 & -1 \\ 0 & 5 & 8 \end{pmatrix}$$

then the augmented matrix

$$\begin{pmatrix} 1 & 2 & 1 & 1 & 0 & 0 \\ -1 & -2 & -2 & 0 & 1 & 0 \\ 0 & 5 & 8 & 0 & 0 & 1 \end{pmatrix}$$

reduces to

$$\begin{pmatrix} 1 & 2 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 5 & 8 & 0 & 0 & 1 \end{pmatrix} \quad (12.2.e)$$

Although this matrix can be row-reduced further, it is not necessary to do so since in equation form we have:

$$\begin{array}{lll} x_{11} + 2x_{21} + x_{31} = 1 & x_{12} + 2x_{22} + x_{32} = 0 & x_{13} + 2x_{23} + x_{33} = 0 \\ \text{(i)} \quad 0 = 1 & \text{(ii)} \quad 0 = 1 & \text{(iii)} \quad 0 = 0 \\ 5x_{21} + 8x_{31} = 0 & 5x_{22} + 8x_{32} = 0 & 5x_{23} + 8x_{33} = 1 \end{array}$$

Clearly, there is no solution to Systems (i) and (ii), therefore A^{-1} does not exist. From this discussion it should be obvious to the reader that the zero row of the coefficient matrix together with the nonzero entry in the fourth column of that row in matrix 12.2e tells us that A^{-1} does not exist.

EXERCISES FOR SECTION 12.2

A Exercises

- In order to develop an understanding of the technique of this section, work out all the details of Example 12.2.1.
- Use the method of this section to find the inverses of the following matrices whenever possible. If an inverse does not exist, explain why.

$$\text{(a)} \begin{pmatrix} 1 & 2 \\ -1 & 3 \end{pmatrix} \quad \text{(b)} \begin{pmatrix} 0 & 3 & 2 & 5 \\ 1 & -1 & 4 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 3 & -1 \end{pmatrix}$$

$$\text{(c)} \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{pmatrix} \quad \text{(d)} \begin{pmatrix} 1 & 2 & 1 \\ -2 & -3 & -1 \\ 1 & 4 & 4 \end{pmatrix}$$

$$\text{(e)} \begin{pmatrix} 6 & 7 & 2 \\ 4 & 2 & 1 \\ 6 & 1 & 1 \end{pmatrix} \quad \text{(f)} \begin{pmatrix} 2 & 1 & 3 \\ 4 & 2 & 1 \\ 8 & 2 & 4 \end{pmatrix}$$

- Same as question 2:

$$\text{(a)} \begin{pmatrix} \frac{1}{3} & 2 \\ \frac{1}{5} & -1 \end{pmatrix} \quad \text{(b)} \begin{pmatrix} 1 & 0 & 0 & 3 \\ 2 & -1 & 0 & 6 \\ 0 & 2 & 1 & 0 \\ 0 & -1 & 3 & 2 \end{pmatrix}$$

$$\text{(c)} \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix} \quad \text{(d)} \begin{pmatrix} 1 & 0 & 0 \\ 2 & 2 & -1 \\ 1 & -1 & 1 \end{pmatrix}$$

$$\text{(e)} \begin{pmatrix} 2 & 3 & 4 \\ 3 & 4 & 5 \\ 4 & 5 & 6 \end{pmatrix} \quad \text{(f)} \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \end{pmatrix}$$

- (a) Find the inverses of the following matrices.

$$\text{(i)} \begin{pmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 5 \end{pmatrix} \quad \text{(ii)} \begin{pmatrix} -1 & 0 & 0 & 0 \\ 0 & \frac{5}{2} & 0 & 0 \\ 0 & 0 & \frac{1}{7} & 0 \\ 0 & 0 & 0 & \frac{3}{4} \end{pmatrix}$$

- If D is a diagonal matrix whose diagonal entries are nonzero, what is D^{-1} ?

5. Express each system of equations in Exercise 1, Section 12.1, in the form $Ax = B$. Solve each system by first finding A^{-1} whenever possible.

12.3 An Introduction to Vector Spaces

When we encountered various types of matrices in Chapter 5, it became apparent that a particular kind of matrix, the diagonal matrix, was much easier to use in computations. For example, if $A = \begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix}$, then A^5 can be found, but its computation is tedious. If

$$D = \begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix}$$

then

$$D^5 = \begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix}^5 = \begin{pmatrix} 1^5 & 0 \\ 0 & 4^5 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1024 \end{pmatrix}$$

In a variety of applications it is beneficial to be able to diagonalize a matrix. In this section we will investigate what this means and consider a few applications. In order to understand when the diagonalization process can be performed, it is necessary to develop several of the underlying concepts of *linear algebra*.

By now, you realize that mathematicians tend to generalize. Once we have found a "good thing," something that is useful, we apply it to as many different concepts as possible. In doing so, we frequently find that the "different concepts" are not really different but only look different. Four sentences in four different languages might look dissimilar, but when they are translated into a common language, they might very well express the exact same idea.

Early in the development of mathematics, the concept of a vector led to a variety of applications in physics and engineering. We can certainly picture vectors, or "arrows," in the xy -plane and even in the three-dimensional space. Does it make sense to talk about vectors in four-dimensional space, in ten-dimensional space, or in any other mathematical situation? If so, what is the essence of a vector? Is it its shape or the rules it follows? The shape in two- or three-space is just a picture, or geometric interpretation, of a vector. The essence is the rules, or properties, we wish vectors to follow so we can manipulate them algebraically. What follows is a definition of what is called a *vector space*. It is a list of all the essential properties of vectors, and it is the basic definition of the branch of mathematics called linear algebra.

Definition: Vector Space. Let V be any nonempty set of objects. Define on V an operation, called addition, for any two elements $\vec{x}, \vec{y} \in V$, and denote this operation by $\vec{x} + \vec{y}$. Let scalar multiplication be defined for a real number $a \in \mathbb{R}$ and any element $\vec{x} \in V$ and denote this operation by $a\vec{x}$. The set V together with operations of addition and scalar multiplication is called a *vector space over \mathbb{R}* if the following hold for all $\vec{x}, \vec{y}, \vec{z} \in V$, and $a, b \in \mathbb{R}$:

- (1) $\vec{x} + \vec{y} = \vec{y} + \vec{x}$
- (2) $(\vec{x} + \vec{y}) + \vec{z} = \vec{x} + (\vec{y} + \vec{z})$
- (3) There exists a vector $\vec{0} \in V$, such that $\vec{x} + \vec{0} = \vec{x}$
- (4) For each vector $\vec{x} \in V$, there exists a unique vector $-\vec{x} \in V$, such that $-\vec{x} + \vec{x} = \vec{0}$.

These are the main properties associated with the operation of addition. They can be summarized by saying that $[V; +]$ is an abelian group.

The next five properties are associated with the operation of scalar multiplication and how it relates to vector addition.

- (5) $a(\vec{x} + \vec{y}) = a\vec{x} + a\vec{y}$
- (6) $(a + b)\vec{x} = a\vec{x} + b\vec{x}$
- (7) $a(b\vec{x}) = (ab)\vec{x}$
- (8) $1\vec{x} = \vec{x}$.

In a vector space it is common to call the elements of V *vectors* and those from \mathbb{R} *scalars*. Vector spaces over the real numbers are also called *real vector spaces*.

Example 12.3.1. Let $V = M_{2 \times 3}(\mathbb{R})$ and let the operations of addition and scalar multiplication be the usual operations of addition and scalar multiplication on matrices. Then V together with these operations is a real vector space. The reader is strongly encouraged to verify the definition for this example before proceeding further (see Exercise 3 of this section). Note we can call the elements of $M_{2 \times 3}(\mathbb{R})$ *vectors* even though they are not arrows.

Example 12.3.2. Let $\mathbb{R}^2 = \{(a_1, a_2) \mid a_1, a_2 \in \mathbb{R}\}$. If we define addition and scalar multiplication the natural way, that is, as we would on 1×2 matrices, then \mathbb{R}^2 is a vector space over \mathbb{R} . (See Exercise 4 of this section.)

In this example, we have the "bonus" that we can illustrate the algebraic concept geometrically. In mathematics, a "geometric bonus" does not always occur and is not necessary for the development or application of the concept. However, geometric illustrations are quite useful in helping us understand concepts and should be utilized whenever available.

Let's consider some illustrations of the vector space \mathbb{R}^2 . Let $\vec{x} = (1, 4)$ and $\vec{y} = (3, 1)$.

We illustrate the vector (a_1, a_2) as a directed line segment, or "arrow," from the point $(0, 0)$ to the point (a_1, a_2) . The vectors \vec{x} and \vec{y} are as pictured in Figure 12.3.1 together with $\vec{x} + \vec{y} = (1, 4) + (3, 1) = (4, 5)$, which also has the geometric representation as pictured in Figure 12.3.1. The vector $2\vec{x} = 2(1, 4) = (2, 8)$ is a vector in the same direction as \vec{x} , but with twice its length.

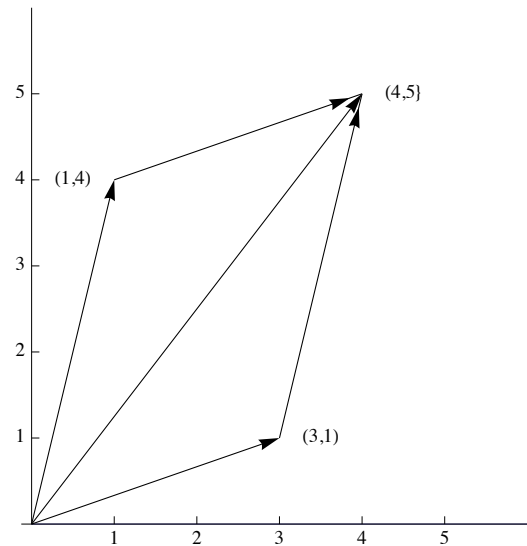


Figure 12.3.1
Addition in \mathbb{R}^2

Remarks:

- (1) We will henceforth drop the arrow above a vector name and use the common convention that boldface letters toward the end of the alphabet are vectors, while letters early in the alphabet are scalars.
- (2) The vector $(a_1, a_2, \dots, a_n) \in \mathbb{R}^n$ is referred to as an n -tuple.
- (3) For those familiar with vector calculus, we are expressing the vector $x = a_1 \mathbf{i} + a_2 \mathbf{j} + a_3 \mathbf{k} \in \mathbb{R}^3$ as (a_1, a_2, a_3) . This allows us to discuss vectors in \mathbb{R}^n in much simpler notation.

In many situations a vector space V is given and we would like to describe the whole vector space by the smallest number of essential reference vectors. An example of this is the description of \mathbb{R}^2 , the xy plane, via the x and y axes. Again our concepts must be algebraic in nature so we are not restricted solely to geometric considerations.

Definition: Linear Combination. A vector \mathbf{y} in vector space V (over \mathbb{R}) is a linear combination of the vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ if there exist scalars a_1, a_2, \dots, a_n in \mathbb{R} such that $\mathbf{y} = a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \dots + a_n \mathbf{x}_n$.

Example 12.3.3 The vector $(2, 3)$ in \mathbb{R}^2 is a linear combination of the vectors $(1, 0)$ and $(0, 1)$ since $(2, 3) = 2(1, 0) + 3(0, 1)$.

Example 12.3.4. Prove that the vector $(5, 4)$ is a linear combination of the vectors $(4, 1)$ and $(1, 3)$. By the definition we must show that there exist scalars a_1 and a_2 such that:

$$(5, 4) = a_1(4, 1) + a_2(1, 3),$$

which reduces to

$$(5, 4) = (4a_1 + a_2, a_1 + 3a_2),$$

which gives us the system of linear equations

$$\begin{aligned} 4a_1 + a_2 &= 5 \\ a_1 + 3a_2 &= 4 \end{aligned}$$

which has solution $a_1 = 1, a_2 = 1$.

Another way of looking at the above example is if we replace a_1 and a_2 both by 1, then the two vectors $(4, 1)$ and $(1, 3)$ produce, or generate, the vector $(5, 4)$. Of course, if we replace a_1 and a_2 by different scalars, we can generate more vectors from \mathbb{R}^2 . If $a_1 = 3$ and $a_2 = -2$, then

$$\begin{aligned} a_1(4, 1) + a_2(1, 3) &= 3(4, 1) + (-2)(1, 3) \\ &= (12, 3) + (-2, -6) \\ &= (12 - 2, 3 - 6) = (10, -3) \end{aligned}$$

Example 12.3.5. Will the vectors $(4, 1)$ and $(1, 3)$ generate any vector we choose in \mathbb{R}^2 ? To see if this is so, we let (b_1, b_2) be an arbitrary vector in \mathbb{R}^2 and see if we can always find scalars a_1 and a_2 such that $a_1(4, 1) + a_2(1, 3) = (b_1, b_2)$. This is equivalent to solving the following system of equations:

$$\begin{aligned} 4a_1 + a_2 &= b_1 \\ a_1 + 3a_2 &= b_2 \end{aligned}$$

which always has solutions for a_1 and a_2 regardless of the values of the real numbers b_1 and b_2 . Why? We formalize in a definition:

Definition: Generate. Let $\{x_1, x_2, \dots, x_n\}$ be a set of vectors in a vector space V over \mathbb{R} . This set is said to generate, or span, V if, for any given vector $y \in V$, we can always find scalars a_1, a_2, \dots, a_n such that $y = a_1 x_1 + a_2 x_2 + \dots + a_n x_n$. A set that generates a vector space is called a generating set.

We now give a geometric interpretation of the above.

We know that the standard coordinate system, x axis and y axis, were introduced in basic algebra in order to describe all points in the xy plane geometrically. It is also quite clear that to describe any point in the plane we need exactly two axes. Form a new coordinate system the following way:

Draw the vector $(4, 1)$ and an axis from the origin through $(4, 1)$ and label it the x' axis. Also draw the vector $(1, 3)$ and an axis from the origin through $(1, 3)$ to be labeled the y' axis. Draw the coordinate grid for the axis, that is, lines parallel, and let the unit lengths of this "new" plane be the lengths of the respective vectors, $(4, 1)$ and $(1, 3)$, so that we obtain Figure 12.3.2.

From Example 12.3.5 and Figure 12.3.2, we see that any vector on the plane can be described using the old (standard xy) axes or our new $x'y'$ axes. Hence the position which had the name $(4, 1)$ in reference to the standard axes has the name $(1, 0)$ with respect to the $x'y'$ axes, or, in the phraseology of linear algebra, the coordinates of the point $(1, 3)$ with respect to the $x'y'$ axes are $(1, 0)$.

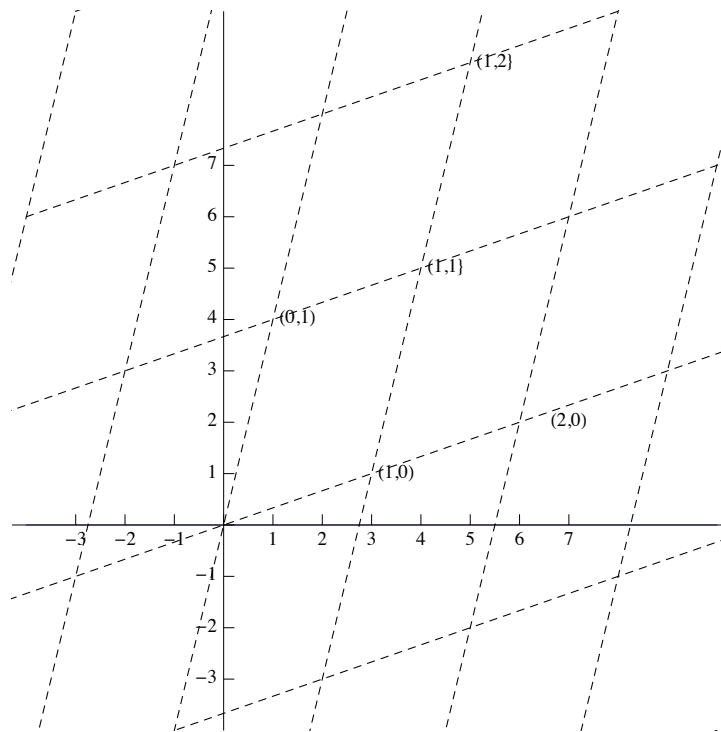


Figure 12.3.2

Example 12.3.6. From Example 12.3.4 we found that if we choose $a_1 = 1$ and $a_2 = 1$, then the two vectors $(4, 1)$ and $(1, 3)$ generate the vector $(5, 4)$. Another geometric interpretation of this problem is that the coordinates of the position $(5, 4)$ with respect to the $x'y'$ axes of Figure 12.3.2 is $(1, 1)$. In other words, a position in the plane has the name $(5, 4)$ in reference to the xy axes and the same position has the name $(1, 1)$ in reference to the $x'y'$ axes.

From the above, it is clear that we can use different axes to describe points or vectors in the plane. No matter what choice we use, we want to be able to describe each position in a unique manner. This is not the case in Figure 12.3.3. Any point in the plane could be described via the $x'y'$ axes, the $x'z'$ axes or the $y'z'$ axes. Therefore, in this case, a single point would have three different names, a very confusing situation.

We formalize the above discussion in two definitions and a theorem.

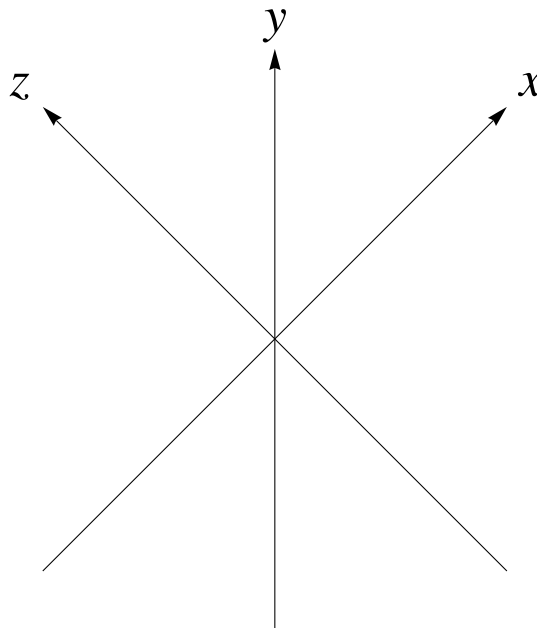


Figure 12.3.3

Definition: Linear Independence/Linear Dependence. The set of vectors $\{x_1, x_2, \dots, x_n\}$ a vector space V (over \mathbb{R}) is linearly independent if the only solution to the equation $a_1 x_1 + a_2 x_2 + \dots + a_n x_n = \mathbf{0}$ is $a_1 = a_2 = \dots = a_n = 0$. Otherwise the set is called a linearly dependent set.

Definition: Basis. A set of vectors $B = \{x_1, x_2, \dots, x_n\}$ is a basis for a vector space V (over \mathbb{R}) if:

- (1) B generates V , and
- (2) B is linearly independent.

Theorem 12.3.1. If $\{x_1, x_2, \dots, x_n\}$ is a basis for a vector space V over \mathbb{R} , then any vector $y \in V$ can be uniquely expressed as a linear combination of the x_i 's.

Proof: Assume that $\{x_1, x_2, \dots, x_n\}$ is a basis for V over \mathbb{R} . We must prove two facts:

- (1) each vector $y \in V$ can be expressed as a linear combination of the x_i 's, and
- (2) each such expression is unique.

Part (1) is trivial since a basis, by its definition, must be a generating set for V .

The proof of (2) is a bit more difficult. We follow the standard approach for any uniqueness facts. Let y be any vector in V and assume that there are two different ways of expressing y , namely

$$y = a_1 x_1 + a_2 x_2 + \dots + a_n x_n$$

and

$$y = b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

where at least one a_i is different from the corresponding b_i . Then equating these two linear combinations we get

$$a_1 x_1 + a_2 x_2 + \dots + a_n x_n = b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

so that

$$(a_1 - b_1) x_1 + (a_2 - b_2) x_2 + \dots + (a_n - b_n) x_n = \mathbf{0}$$

Now a crucial observation: since the x_i 's form a linearly independent set, the only solution to the previous equation is that each of the coefficients must equal zero, so $a_i - b_i = 0$ for $i = 1, 2, \dots, n$. Hence $a_i = b_i$, for all i . This contradicts our assumption that at least one a_i is different from the corresponding b_i , so each vector $y \in V$ can be expressed in one and only one way. ■

Theorem 12.3.1, together with the previous examples, gives us a clear insight into the meaning of linear independence, namely uniqueness.

Example 12.3.7. Prove that $\{(1, 1), (-1, 1)\}$ is a basis for \mathbb{R}^2 over \mathbb{R} and explain what this means geometrically. First we must show that the vectors $(1, 1)$ and $(-1, 1)$ generate all of \mathbb{R}^2 . This we can do by imitating Example 12.3.5 and leave it to the reader (see Exercise 10 of this section). Secondly, we must prove that the set is linearly independent.

Let a_1 and a_2 be scalars such that $a_1(1, 1) + a_2(-1, 1) = (0, 0)$. We must prove that the only solution to the equation is that a_1 and a_2 must both equal zero. The above equation becomes $(a_1 - a_2, a_1 + a_2) = (0, 0)$ which gives us the system

$$\begin{aligned} a_1 - a_2 &= 0 \\ a_1 + a_2 &= 0 \end{aligned}$$

The augmented matrix of this system reduces in such way that the only solution is the trivial one of all zeros:

$$\begin{pmatrix} 1 & -1 & 0 \\ 1 & 1 & 0 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} \Rightarrow a_1 = a_2 = 0$$

Therefore, the set is linearly independent.

To explain the results geometrically, note through Exercise 12, part a, that the coordinates of each vector $\mathbf{y} \in \mathbb{R}^2$ can be determined uniquely using the vectors $(1,1)$ and $(-1, 1)$. The concept of dimension is quite obvious for those vector spaces that have an immediate geometric interpretation. For example, the dimension of \mathbb{R}^2 is two and that of \mathbb{R}^3 is three. How can we define the concept of dimension algebraically so that the resulting definition correlates with that of \mathbb{R}^2 and \mathbb{R}^3 ? First we need a theorem, which we will state without proof.

Theorem 12.3.2. If V is a vector space with a basis containing n elements, then all bases of V contain n elements.

Definition: Dimension. Let V be a vector space over \mathbb{R} with basis $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$. Then the dimension of V is n . We use the notation $\dim V = n$ to indicate that V is n -dimensional

EXERCISES FOR SECTION 12.3

A Exercises

1. If $a = 2, b = -3$,

$$A = \begin{pmatrix} 1 & 0 & -1 \\ 2 & 3 & 4 \end{pmatrix}, \quad B = \begin{pmatrix} 2 & -2 & 3 \\ 4 & 5 & 8 \end{pmatrix}, \quad \text{and } C = \begin{pmatrix} 1 & 0 & 0 \\ 3 & 2 & -2 \end{pmatrix}$$

verify that all properties of the definition of a vector space are true for $M_{2 \times 3}(\mathbb{R})$ with these values.

2. Let $a = 3, b = 4, \mathbf{x} = (-1, 3), \mathbf{y} = (2, 3)$, and $\mathbf{z} = (1, 0)$. Verify that all properties of the definition of a vector space are true for \mathbb{R}^2 for these values.

3. (a) Verify that $M_{2 \times 3}(\mathbb{R})$ is a vector space over \mathbb{R} .

(b) Is $M_{m \times n}(\mathbb{R})$ a vector space over \mathbb{R} ?

4. (a) Verify that \mathbb{R}^2 is a vector space over \mathbb{R} .

(b) Is \mathbb{R}^n a vector space over \mathbb{R} for every positive integer n ?

5. Let $P^3 = \{a_0 + a_1x + a_2x^2 + a_3x^3 \mid a_0, a_1, a_2, a_3 \in \mathbb{R}\}$; that is, P^3 is the set of all polynomials in x having real coefficients with degree less than or equal to 3. Verify that P^3 is a vector space over \mathbb{R} .

6. For each of the following, express the vector \mathbf{y} as a linear combination of the vectors \mathbf{x}_1 and \mathbf{x}_2 .

(a) $\mathbf{y} = (5, 6), \mathbf{x}_1 = (1, 0)$, and $\mathbf{x}_2 = (0, 1)$

(b) $\mathbf{y} = (2, 1), \mathbf{x}_1 = (2, 1)$, and $\mathbf{x}_2 = (1, 1)$

(c) $\mathbf{y} = (3, 4), \mathbf{x}_1 = (1, 1)$, and $\mathbf{x}_2 = (-1, 1)$

7. Express the vector $\begin{pmatrix} 1 & 2 \\ -3 & 3 \end{pmatrix} \in M_{2 \times 2}(\mathbb{R})$, as a linear combination of

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, \begin{pmatrix} -1 & 5 \\ 2 & 1 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \text{ and } \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$$

8. Express the vector $x^3 - 4x^2 + 3 \in P^3$ as a linear combination of the vectors $1, x, x^2$, and x^3 .

9. (a) Show that the set $\{\mathbf{x}_1, \mathbf{x}_2\}$ generates \mathbb{R}^2 for each of the parts in Exercise 6 of this section.

(b) Show that $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ generates \mathbb{R}^2 where $\mathbf{x}_1 = (1, 1), \mathbf{x}_2 = (3, 4)$, and $\mathbf{x}_3 = (-1, 5)$.

(c) Create a set of four or more vectors that generates \mathbb{R}^2 .

(d) What is the smallest number of vectors needed to generate \mathbb{R}^2 ? \mathbb{R}^n ?

(e) Show that the set of matrices containing

$$\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} \text{ and } \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$$

generates $M_{2 \times 2}(\mathbb{R})$

(f) Show that $\{1, x, x^2, x^3\}$ generates P^3 .

10. Complete Example 12.3.7 by showing that $\{(1, 1), (-1, 1)\}$ generates \mathbb{R}^2

11. (a) Prove that $\{(4, 1), (1, 3)\}$ is a basis for \mathbb{R}^2 over \mathbb{R} .

(b) Prove that $\{(1, 0), (3, 4)\}$ is a basis for \mathbb{R}^2 over \mathbb{R} .

(c) Prove that $\{(1, 0, -1), (2, 1, 1), (1, -3, -1)\}$ is a basis for \mathbb{R}^3 over \mathbb{R} .

(d) Prove that the sets in Exercise 9, parts e and f, form bases of the respective vector spaces.

12. (a) Determine the coordinates of the points or vectors $(3, 4)$, $(-1, 1)$, and $(1, 1)$ with respect to the basis $\{(1, 1), (-1, 1)\}$ of \mathbb{R}^2 . Interpret your results geometrically,

(b) Determine the coordinates of the points or vector $(3, 5, 6)$ with respect to the basis $\{(1, 0, 0), (0, 1, 0), (0, 0, 1)\}$. Explain why this basis is called the standard basis for \mathbb{R}^3 .

13. (a) Let $y_1 = (1, 3, 5, 9)$, $y_2 = (5, 7, 6, 3)$, and $c = 2$. Find $y_1 + y_2$ and $c y_1$.

(b) Let $f_1(x) = 1 + 3x + 5x^2 + 9x^3$, $f_2(x) = 5 + 7x + 6x^2 + 3x^3$ and $c = 2$. Find $f_1(x) + f_2(x)$ and $c f_1(x)$.

(c) Let $A = \begin{pmatrix} 1 & 3 \\ 5 & 9 \end{pmatrix}$, $B = \begin{pmatrix} 5 & 7 \\ 6 & 3 \end{pmatrix}$, and $c = 2$. Find $A + B$ and $c A$.

(d) Are the vector spaces \mathbb{R}^4 , P^3 and $M_{2 \times 2}(\mathbb{R})$ isomorphic to each other? Discuss with reference to parts a, b, and c.

12.4 The Diagonalization Process

We now have the background to understand the main ideas behind the diagonalization process.

Definition: Eigenvalue, Eigenvector. Let A be an $n \times n$ matrix over \mathbb{R} . λ is an eigenvalue of A if for some nonzero column vector $\mathbf{x} \in \mathbb{R}^n$ we have $A\mathbf{x} = \lambda\mathbf{x}$. \mathbf{x} is called an Eigenvector corresponding to the eigenvalue λ .

Example 12.4.1. Find the eigenvalues and corresponding eigenvectors of the matrix $A = \begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix}$. We want to find nonzero vectors $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ and real numbers λ such that

$$\begin{aligned} A\mathbf{x} = \lambda\mathbf{x} &\Leftrightarrow \begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \lambda \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \\ &\Leftrightarrow \begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \lambda \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ &\Leftrightarrow \begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ &\Leftrightarrow \left(\begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ &\Leftrightarrow \begin{pmatrix} 2-\lambda & 1 \\ 2 & 3-\lambda \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad (12.4a) \end{aligned}$$

The last matrix equation will have nonzero solutions if and only if

$$\det \begin{pmatrix} 2-\lambda & 1 \\ 2 & 3-\lambda \end{pmatrix} = 0$$

or $(2-\lambda)(3-\lambda) - 2 = 0$, which simplifies to $\lambda^2 - 5\lambda + 4 = 0$. Therefore, the solutions to this quadratic equation, $\lambda_1 = 1$ and $\lambda_2 = 4$, are the eigenvalues of A . We now have to find eigenvectors associated with each eigenvalue.

Case 1. For $\lambda_1 = 1$, Equation 12.4a becomes:

$$\begin{aligned} \begin{pmatrix} 2-1 & 1 \\ 2 & 3-1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ \begin{pmatrix} 1 & 1 \\ 2 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \end{aligned}$$

which reduces to the single equation, $x_1 + x_2 = 0$. From this, $x_1 = -x_2$. This means the solution set of this equation is (in column notation)

$$E_1 = \left\{ \begin{pmatrix} -c \\ c \end{pmatrix} \mid c \in \mathbb{R} \right\}$$

So any column vector of the form $\begin{pmatrix} -c \\ c \end{pmatrix}$ where c is any nonzero real number is an eigenvector associated with $\lambda_1 = 1$. The reader should verify that, for example,

$$\begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix} \begin{pmatrix} \frac{2}{3} \\ -\frac{2}{3} \end{pmatrix} = 1 \begin{pmatrix} \frac{2}{3} \\ -\frac{2}{3} \end{pmatrix}$$

so that $\begin{pmatrix} \frac{2}{3} \\ -\frac{2}{3} \end{pmatrix}$ is an eigenvector associated with eigenvalue 1.

Case 2. For $\lambda_2 = 4$ equation 12.4.a becomes:

$$\begin{aligned} \begin{pmatrix} 2-4 & 1 \\ 2 & 3-4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ \begin{pmatrix} -2 & 1 \\ 2 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} &= \begin{pmatrix} 0 \\ 0 \end{pmatrix} \end{aligned}$$

which reduces to the single equation $-2x_1 + x_2 = 0$, so that $x_2 = 2x_1$. The solution set of the equation is

$$E_2 = \left\{ \begin{pmatrix} c \\ 2c \end{pmatrix} \mid c \in \mathbb{R} \right\}$$

Therefore, all eigenvectors of A associated with the eigenvalue $\lambda_2 = 4$ are of the form $\begin{pmatrix} c \\ -2c \end{pmatrix}$, where c can be any nonzero number.

The following theorems summarize the more important aspects of this example:

Theorem 12.4.1. Let A be any $n \times n$ matrix over \mathbb{R} . Then $\lambda \in \mathbb{R}$ is an eigenvalue of A if and only if $\det(A - \lambda I) = 0$.

The equation $\det(A - \lambda I) = 0$ is called the *characteristic equation* and the left side of this equation is called the *characteristic polynomial* of A .

Theorem 12.4.2. Nonzero eigenvectors corresponding to distinct eigenvalues are linearly independent.

The solution space of $(A - \lambda I)\mathbf{x} = \mathbf{0}$ is called the *eigenspace of A corresponding to λ* . This terminology is justified by Exercise 2 of this section.

We now consider the main aim of this section. Given an $n \times n$ (square) matrix A , we would like to "change" A into a diagonal matrix D , perform our tasks with the simpler matrix D , and then describe the results in terms of the given matrix A .

Definition: Diagonalizable Matrix. An $n \times n$ matrix A is called *diagonalizable* if there exists an invertible $n \times n$ matrix P such that $P^{-1}AP$ is a diagonal matrix D . The matrix P is said to *diagonalize the matrix A* .

Example 12.4.2. We will now diagonalize the matrix A of Example 12.4.1. Form the matrix P as follows: Let $P^{(1)}$ be the first column of P . Choose for $P^{(1)}$ any eigenvector from E_1 . We may as well choose a simple vector in E_1 so $P^{(1)} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ is our candidate. Similarly, let $P^{(2)}$ be the second

column of P , and choose for $P^{(2)}$ any eigenvector from E_2 . The vector $P^{(2)} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ is a reasonable choice, thus

$$P = \begin{pmatrix} 1 & 1 \\ -1 & 2 \end{pmatrix} \quad \text{and} \quad P^{-1} = \frac{1}{3} \begin{pmatrix} 2 & -1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} \frac{2}{3} & -\frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} \end{pmatrix}$$

So that

$$P^{-1}AP = \frac{1}{3} \begin{pmatrix} 2 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ -1 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix}$$

Notice that the elements on the main diagonal of D are the eigenvalues of A , where D_{ii} is the eigenvalue corresponding to the eigenvector $P^{(i)}$.

Remarks:

(1) The first step in the diagonalization process is the determination of the eigenvalues. The ordering of the eigenvalues is purely arbitrary. If we designate $\lambda_1 = 4$ and $\lambda_2 = 1$, the columns of P would be interchanged and D would be $\begin{pmatrix} 4 & 0 \\ 0 & 1 \end{pmatrix}$ (see Exercise 3b of this section). Nonetheless, the final outcome of the application to which we are applying the diagonalization process would be the same.

(2) If A is an $n \times n$ matrix with distinct eigenvalues, then P is also an $n \times n$ matrix whose columns $P^{(1)}, P^{(2)}, \dots, P^{(n)}$ are n linearly independent vectors.

Example 12.4.3. Diagonalize the matrix

$$A = \begin{pmatrix} 1 & 12 & -18 \\ 0 & -11 & 18 \\ 0 & -6 & 10 \end{pmatrix}.$$

$$\begin{aligned} \det(A - \lambda I) &= \det \begin{pmatrix} 1-\lambda & 12 & -18 \\ 0 & -\lambda-11 & 18 \\ 0 & -6 & 10-\lambda \end{pmatrix} \\ &= (1-\lambda) \det \begin{pmatrix} -\lambda-11 & 18 \\ -6 & 10-\lambda \end{pmatrix} \\ &= (1-\lambda)((-\lambda-11)(10-\lambda) + 108) \\ &= (1-\lambda)(\lambda^2 + \lambda - 2) \end{aligned}$$

Hence, the equation $\det(A - \lambda I) = 0$ becomes

$$(1-\lambda)(\lambda^2 + \lambda - 2) = -(\lambda-1)^2(\lambda+2)$$

Therefore, our eigenvalues for A are $\lambda_1 = -2$ and $\lambda_2 = 1$. We note that we do not have three distinct eigenvalues, but we proceed as in the previous example.

Case 1. For $\lambda_1 = -2$ the equation $(A - \lambda I)\mathbf{x} = \mathbf{0}$ becomes

$$\begin{pmatrix} 3 & 12 & -18 \\ 0 & -9 & 18 \\ 0 & -6 & 12 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

Using *Mathematica*, we can row reduce the matrix:

$$\text{RowReduce} \left[\begin{pmatrix} 3 & 12 & -18 \\ 0 & -9 & 18 \\ 0 & -6 & 12 \end{pmatrix} \right]$$

$$\begin{pmatrix} 1 & 0 & 2 \\ 0 & 1 & -2 \\ 0 & 0 & 0 \end{pmatrix}$$

In equation form, the matrix equation is then equivalent to

$$\begin{aligned} x_1 &= -2x_3 \\ x_2 &= 2x_3 \end{aligned}$$

Therefore, the solution, or eigenspace, corresponding to $\lambda_1 = -2$ consists of vectors of the form

$$\begin{pmatrix} -2x_3 \\ 2x_3 \\ x_3 \end{pmatrix} = x_3 \begin{pmatrix} -2 \\ 2 \\ 1 \end{pmatrix}$$

Therefore $\begin{pmatrix} -2 \\ 2 \\ 1 \end{pmatrix}$ is an eigenvector corresponding to the eigenvalue $\lambda_1 = -2$, and can be used for our first column of P :

$$P^{(1)} = \begin{pmatrix} -2 \\ 2 \\ 1 \end{pmatrix}$$

Before we continue we make the observation: E_2 is a subspace of \mathbb{R}^3 with basis $\{P^{(1)}\}$ and $\dim E_1 = 1$.

Case 2. If $\lambda_2 = 1$, then the equation $(A - \lambda I)x = \mathbf{0}$ becomes

$$\begin{pmatrix} 0 & 12 & -18 \\ 0 & -12 & 18 \\ 0 & -6 & 9 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

Without the aid of any computer technology, it should be clear that all three equations that correspond to this matrix equation are equivalent to $2x_2 - 3x_3 = 0$, or $x_2 = \frac{3}{2}x_3$. Notice that x_1 can take on any value, so any vector of the form

$$\begin{pmatrix} x_1 \\ \frac{3}{2}x_3 \\ x_3 \end{pmatrix} = x_1 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + x_3 \begin{pmatrix} 0 \\ \frac{3}{2} \\ 1 \end{pmatrix}$$

will solve the matrix equation.

We note that the solution set contains two independent variables, x_1 and x_3 . Further, note that we cannot express the eigenspace E_2 as a linear combination of a single vector as in Case 1. However, it can be written as

$$E_2 = \left\{ x_1 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + x_3 \begin{pmatrix} 0 \\ \frac{3}{2} \\ 1 \end{pmatrix} \mid x_1, x_3 \in \mathbb{R} \right\}.$$

We can replace any vector in a basis with a nonzero multiple of that vector. Simply for aesthetic reasons, we will multiply the second vector that generates E_2 by 2. Therefore, the eigenspace E_2 is a subspace of \mathbb{R}^3 with basis $\left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 3 \\ 2 \end{pmatrix} \right\}$ and so $\dim E_2 = 2$.

What this means with respect to the diagonalization process is that $\lambda_2 = 1$ gives us both Column 2 and Column 3 the diagonalizing matrix. The order is not important. Let

$$P^{(2)} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \text{ and } P^{(3)} = \begin{pmatrix} 0 \\ 3 \\ 2 \end{pmatrix} \text{ and so } P = \begin{pmatrix} -2 & 1 & 0 \\ 2 & 0 & 3 \\ 1 & 0 & 2 \end{pmatrix}$$

The reader can verify (see Exercise 5 of this section) that

$$P^{-1} = \begin{pmatrix} 0 & 2 & -3 \\ 1 & 4 & -6 \\ 0 & -1 & 2 \end{pmatrix} \text{ and } P^{-1}AP = \begin{pmatrix} -2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

In doing Example 12.4.3, the given 3×3 matrix A produced only two, not three, distinct eigenvalues, yet we were still able to diagonalize A . The reason we were able to do so was because we were able to find three linearly independent eigenvectors. Again, the main idea is to produce a matrix P that does the diagonalizing. If A is an $n \times n$ matrix, P will be an $n \times n$ matrix, and its n columns must be linearly independent eigenvectors. The main question in the study of diagonalizability is “When can it be done?” This is summarized in the following theorem.

Theorem 12.4.3. *Let A be an $n \times n$ matrix. Then A is diagonalizable if and only if A has n linearly independent eigenvectors.*

Outline of a proof: (\Leftarrow) Assume that A has linearly independent eigenvectors, $P^{(1)}, P^{(2)}, \dots, P^{(n)}$, with corresponding eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_n$. We want to prove that A is diagonalizable. Column i of the $n \times n$ matrix AP is $AP^{(i)}$ (see Exercise 7 of this section). Then, since the $P^{(i)}$ is an eigenvector of A associated with the eigenvalue λ_i we have $AP^{(i)} = \lambda_i P^{(i)}$ for $i = 1, 2, \dots, n$. But this means that $AP = PD$, where D is the diagonal matrix with diagonal entries $\lambda_1, \lambda_2, \dots, \lambda_n$. If we multiply both sides of the equation by P^{-1} we get the desired $P^{-1}AP = D$.

(\Rightarrow) The proof in this direction involves a concept that is not covered in this text (rank of a matrix); so we refer the interested reader to virtually any linear algebra text for a proof. ■

We now give an example of a matrix which is not diagonalizable.

Example 12.4.4. Let us attempt to diagonalize the matrix $A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 1 \\ 1 & -1 & 4 \end{pmatrix}$

$$A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 1 \\ 1 & -1 & 4 \end{pmatrix}$$

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 1 \\ 1 & -1 & 4 \end{pmatrix}$$

$$\begin{aligned} \det(A - \lambda I) &= \det \begin{pmatrix} 1-\lambda & 0 & 0 \\ 0 & 2-\lambda & 1 \\ 1 & -1 & 4-\lambda \end{pmatrix} \\ &= (1-\lambda) \det \begin{pmatrix} 2-\lambda & 1 \\ -1 & 4-\lambda \end{pmatrix} \\ &= (1-\lambda)((2-\lambda)(4-\lambda) + 1) \\ &= (1-\lambda)(\lambda^2 - 6\lambda + 9) \\ &= (1-\lambda)(\lambda - 3)^2 \end{aligned}$$

$$\det(A - \lambda I) = 0 \Rightarrow \lambda = 1 \text{ or } \lambda = 3$$

Therefore there are two eigenvalues, $\lambda_1 = 1$ and $\lambda_2 = 3$. Since λ_1 is an eigenvalue of degree 1 it will have an eigenspace of dimension 1. Since λ_2 is a double root of the characteristic equation, the dimension of its eigenspace must be 2 in order to be able to diagonalize.

Case 1. For $\lambda_1 = 1$, the equation $(A - \lambda I)\mathbf{x} = \mathbf{0}$ becomes

$$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 1 & -1 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

A quick *Mathematica* evaluation make the solution to this system obvious

RowReduce[A - IdentityMatrix[3]]

$$\begin{pmatrix} 1 & 0 & 4 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix}$$

There is one free variable, x_3 , and

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} -4x_3 \\ -x_3 \\ x_3 \end{pmatrix} = x_3 \begin{pmatrix} -4 \\ -1 \\ 1 \end{pmatrix}$$

Hence, $\left\{ \begin{pmatrix} -4 \\ -1 \\ 1 \end{pmatrix} \right\}$ is a basis for the eigenspace of $\lambda_1 = 1$.

Case 2. For $\lambda_2 = 3$, the equation $(A - \lambda I)\mathbf{x} = \mathbf{0}$ becomes

$$\begin{pmatrix} -2 & 0 & 0 \\ 0 & -1 & 1 \\ 1 & -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

RowReduce[A - 3 IdentityMatrix[3]]

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{pmatrix}$$

Once again there is only one free variable in the row reduction and so the dimension of the eigenspace will be one:

$$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ x_3 \\ x_3 \end{pmatrix} = x_3 \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$$

Hence, $\left\{ \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix} \right\}$ is a basis for the eigenspace of $\lambda_2 = 3$. This means that $\lambda_2 = 3$ produces only one column for P . Since we began with only two eigenvalues, we had hoped that one of them would produce a vector space of dimension two, or, in matrix terms, two linearly independent columns of P . Since A does not have three linearly independent eigenvectors A cannot be diagonalized.



Mathematica Note

Diagonalization can be easily done with a few built-in functions of *Mathematica*. Here is a 3×3 matrix we've selected because the eigenvalues are very simple, and could be found by hand with a little work.

$$\mathbf{A} = \begin{pmatrix} 4 & 1 & 0 \\ 1 & 5 & 1 \\ 0 & 1 & 4 \end{pmatrix};$$

The set of linearly independent eigenvectors of A can be computed:

Eigenvectors[A]

$$\begin{pmatrix} 1 & 2 & 1 \\ -1 & 0 & 1 \\ 1 & -1 & 1 \end{pmatrix}$$

The rows of this matrix are the eigenvectors, so we transpose the result to get our diagonalizing matrix P whose columns are eigenvectors.

P = Transpose[Eigenvectors[A]]

$$\begin{pmatrix} 1 & -1 & 1 \\ 2 & 0 & -1 \\ 1 & 1 & 1 \end{pmatrix}$$

We then use P to diagonalize. The entries in the diagonal matrix are the eigenvalues of A .

Inverse[P].A.P

$$\begin{pmatrix} 6 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 3 \end{pmatrix}$$

We could have gotten the eigenvalues directly this way:

Eigenvalues[A]

{6, 4, 3}

Most matrices that are selected at random will not have "nice" eigenvalues. Here is a new matrix A that looks similar to the one above.

$$A = \begin{pmatrix} 8 & 1 & 0 \\ 1 & 5 & 1 \\ 0 & 1 & 7 \end{pmatrix};$$

Asking for the eigenvalues first, we see that the result is returned symbolically as the three roots to a cubic equation. The default for *Mathematica* is to leave these uncomputed. Since the entries of A are exact numbers, *Mathematica* is capable of giving an exact solution, but it's very messy. The easiest way around the problem is to make the entries in A approximate. The following expression redefines A as approximate.

A = N[A]

$$\begin{pmatrix} 8. & 1. & 0. \\ 1. & 5. & 1. \\ 0. & 1. & 7. \end{pmatrix}$$

Now we can get approximate eigenvalues, and the approximations are very good for most purposes.

Eigenvalues[A]

{8.3772, 7.27389, 4.34891}

We can verify that the matrix can be diagonalized although due to round-off error some of the off-diagonal entries of the "diagonal" matrix are nonzero.

P = Transpose[Eigenvectors[A]]

$$\begin{pmatrix} 0.906362 & -0.341882 & 0.248244 \\ 0.341882 & 0.248244 & -0.906362 \\ 0.248244 & 0.906362 & 0.341882 \end{pmatrix}$$

Inverse[P].A.P

$$\begin{pmatrix} 8.3772 & 2.22045 \times 10^{-16} & 6.66134 \times 10^{-16} \\ 0. & 7.27389 & 4.44089 \times 10^{-16} \\ 1.66533 \times 10^{-15} & -4.44089 \times 10^{-16} & 4.34891 \end{pmatrix}$$

The **Chop** function will set small numbers to zero. The default threshold for "small" is 10^{-10} but that can be adjusted, if desired.

Diag = Chop[Inverse[P].A.P]

$$\begin{pmatrix} 8.3772 & 0 & 0 \\ 0 & 7.27389 & 0 \\ 0 & 0 & 4.34891 \end{pmatrix}$$

We can't use the name **D** here because *Mathematica* reserves it for the differentiation function.

If you experiment with more matrices, you will undoubtedly encounter situations where some eigenvalues are complex. The process is the same, although we've avoided these just for simplicity.



Sage Note

We start by defining the same matrix as we did in *Mathematica*. We also declare D and P to be variables.

```
A = Matrix(QQ, [[4, 1, 0], [1, 5, 1], [0, 1, 4]]); A
[4 1 0]
[1 5 1]
[0 1 4]
var ('D, P')
(D, P)
```

We have been working with "right eigenvectors" since the x in $Ax = \lambda x$ is a column vector to the right of A . It's not so common but still desirable in some situations to consider "left eigenvectors," so Sage allows either one. The `right_eigenmatrix` method returns a pair of matrices. The diagonal matrix, D , with eigenvalues and the diagonalizing matrix, P , which is made up of columns that are eigenvectors corresponding to the eigenvectors of D .

```
(D,P)=A.right_eigenmatrix(); (D,P)
```

$$\begin{pmatrix} 6 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{pmatrix}$$

```
[0 4 0] [ 2 0 -1]
[0 0 3], [ 1 -1 1]
)
```

We should note here that \mathbf{P} is not unique because even if an eigenspace has dimension one, any nonzero vector in that space will serve as an eigenvector. For that reason, the \mathbf{P} generated by Sage isn't identical to the one generated by *Mathematica*, but they both work. Here we verify the result for our Sage calculation. Recall that an asterisk is used for matrix multiplication in Sage.

```
P.inverse()*A*P
=
[6 0 0]
[0 4 0]
[0 0 3]
```

Here is a second matrix, again the same as we used with *Mathematica*.

```
A2=Matrix(QQ,[[8,1,0],[1,5,1],[0,1,7]]);A2
[8 1 0]
[1 5 1]
[0 1 7]
```

Here we've already specified that the underlying system is the rational numbers. Since the eigenvalues are not rational, Sage will revert to approximate number by default. We'll just pull out the matrix of eigenvectors this time and display rounded entries. Here the diagonalizing matrix looks very different from the result from *Mathematica*, but this is because the eigenvalues are not in the same order in the two calculations. They both diagonalize but with a different diagonal matrix.

```
P=A2.right_eigenmatrix()[1]
P.numerical_approx(digits=3)
[ 1.00 1.00 1.00]
[ -3.65 -0.726 0.377]
[ 1.38 -2.65 0.274]
D=(P.inverse()*A2*P);D.numerical_approx(digits=3)
[ 4.35 0.000 0.000]
[0.000 7.27 0.000]
[0.000 0.000 8.38]
```

EXERCISES FOR SECTION 12.4

A Exercises

1. (a) List three different eigenvectors of $A = \begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix}$, the matrix of Example 12.4.1, associated with the two eigenvalues 1 and 4. Verify your results.

(b) Choose one of the three eigenvectors corresponding to 1 and one of the three eigenvectors corresponding to 4, and show that the two chosen vectors are linearly independent.

2. (a) Verify that E_1 and E_2 in Example 12.4.1 are vector spaces over \mathbb{R} . Since they are also subsets of \mathbb{R}^2 , they are called subvector-spaces, or subspaces for short, of \mathbb{R}^2 . Since these are subspaces consisting of eigenvectors, they are called eigenspaces.

(b) Use the definition of dimension in the previous section to find $\dim E_1$ and $\dim E_2$. Note that $\dim E_1 + \dim E_2 = \dim \mathbb{R}^2$. This is not a coincidence.

3. (a) Verify that $P^{-1}AP$ is indeed equal to $\begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix}$, as indicated in Example 12.4.2.

(b) Choose $P^{(1)} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ and $P^{(2)} = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ and verify that the new value of P satisfies $P^{-1}AP = \begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix}$.

(c) Take any two linearly independent eigenvectors of the matrix A of Example 12.4.2 and verify that $P^{-1}AP$ is a diagonal matrix.

4. (a) Let A be the matrix in Example 12.4.3 and $P = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 2 \end{pmatrix}$. Without doing any actual matrix multiplications, determine the value of

$P^{-1}AP$

(b) If you choose the columns of P in the reverse order, what is $P^{-1}AP$?

5. Diagonalize the following, if possible:

(a) $\begin{pmatrix} 1 & 2 \\ 3 & 2 \end{pmatrix}$ (b) $\begin{pmatrix} -2 & 1 \\ -7 & 6 \end{pmatrix}$ (c) $\begin{pmatrix} 3 & 0 \\ 0 & 4 \end{pmatrix}$

$$(d) \begin{pmatrix} 1 & -1 & 4 \\ 3 & 2 & -1 \\ 2 & 1 & -1 \end{pmatrix} \quad (e) \begin{pmatrix} 6 & 0 & 0 \\ 0 & 7 & -4 \\ 9 & 1 & 3 \end{pmatrix} \quad (f) \begin{pmatrix} 1 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 1 \end{pmatrix}$$

6. Diagonalize the following, if possible:

$$(a) \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \quad (b) \begin{pmatrix} 2 & 1 \\ 4 & 2 \end{pmatrix} \quad (c) \begin{pmatrix} 2 & -1 \\ 1 & 0 \end{pmatrix}$$

$$(d) \begin{pmatrix} 1 & 3 & 6 \\ -3 & -5 & -6 \\ 3 & 3 & 6 \end{pmatrix} \quad (e) \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix} \quad (f) \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{pmatrix}$$

B Exercise

7. Let A and P be as in Example 12.4.3. Show that the columns of the matrix AP can be found by computing $AP^{(1)}, AP^{(2)}, \dots, AP^{(n)}$.

8. Prove that if P is an $n \times n$ matrix and D is a diagonal matrix with diagonal entries d_1, d_2, \dots, d_n , then PD is the matrix obtained from P , but multiplying column i of P by $d_i, i = 1, 2, \dots, n$.

C Exercise

9. (a) There is an option to the *Mathematica* functions **Eigenvectors** and **Eigenvalues** called **Cubics** that will use the cubic equation to find exact eigenvalues of a matrix like $\begin{pmatrix} 8 & 1 & 0 \\ 1 & 5 & 1 \\ 0 & 1 & 7 \end{pmatrix}$. Use that option to find the exact eigenvalues of the matrix. Diagonalize the matrix using the **Cubics** option and then convert the result to a matrix of approximate numbers to compare your result with the approximate result we found in the *Mathematica* Note.

12.5 Some Applications

A large and varied number of applications involve computations of powers of matrices. These applications can be found in science, the social sciences, economics, the analysis of relationships with groups, engineering, and, indeed, any area where mathematics is used and, therefore, where programs are to be developed. We will consider a few diverse examples here.

To aid your understanding of the following examples, we develop a helpful technique to compute $A^m, m > 1$. If A can be diagonalized, then there is a matrix P such that $P^{-1}AP = D$, where D is a diagonal matrix and

$$A^m = P D^m P^{-1} \text{ for all } m \geq 1. \quad (12.5 a)$$

You are asked to prove this equation in Exercise 9 of Section 5.4. The condition that D be a diagonal matrix is not necessary but when it is, the calculation on the right side is particularly easy to perform. Although the formal proof of equation 12.4a is done by induction, the reason *why* it is true is easily seen by writing out an example such as $m = 3$:

$$\begin{aligned} A^m &= (P D P^{-1})^m \quad \text{To get this, solve } P^{-1}AP = D \text{ for } A \text{ and substitute} \\ &= (P D P^{-1})(P D P^{-1})(P D P^{-1}) \\ &= P D (P^{-1}P) D (P^{-1}P) D P^{-1} \quad \text{by associativity of matrix mult.} \\ &= P D I D I D P^{-1} \\ &= P D D D P^{-1} \\ &= P D^3 P^{-1} \end{aligned}$$

Example 12.5.1: Recursion. Consider the computation of terms of the Fibonacci sequence, which we examined in Example 8.1.5:

$$F_0 = 1, F_1 = 1$$

$$F_k = F_{k-1} + F_{k-2} \text{ for } k \geq 2.$$

In order to formulate the calculation in matrix form, we introduced the "dummy equation" $F_{k-1} = F_{k-1}$ so that now we have two equations

$$\begin{aligned} F_k &= F_{k-1} + F_{k-2} \\ F_{k-1} &= F_{k-1} \end{aligned}$$

These two equations can be expressed in matrix form as

$$\begin{aligned}
\begin{pmatrix} F_k \\ F_{k-1} \end{pmatrix} &= \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} F_{k-1} \\ F_{k-2} \end{pmatrix} \text{ if } k \geq 2 \\
&= A \begin{pmatrix} F_{k-1} \\ F_{k-2} \end{pmatrix} \text{ if } A = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} \\
&= A^2 \begin{pmatrix} F_{k-2} \\ F_{k-3} \end{pmatrix} \text{ if } k \geq 3 \\
&\text{etc. if } k \text{ is large enough}
\end{aligned}$$

We can use induction to prove that if $k \geq 2$,

$$\begin{pmatrix} F_k \\ F_{k-1} \end{pmatrix} = A^{k-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Next, by diagonalizing A and using the fact that $A^m = P D^m P^{-1}$, we can show that

$$F_k = \frac{1}{\sqrt{5}} \left(\left(\frac{1+\sqrt{5}}{2} \right)^k - \left(\frac{1-\sqrt{5}}{2} \right)^k \right)$$

See Exercise 1a of this section.

Comments:

(1) An equation of the form $F_k = aF_{k-1} + bF_{k-2}$, where a and b are given constants, is referred to linear homogeneous second-order difference equation. The conditions $F_0 = c_0$ and $F_1 = c_1$, where c_1 and c_2 are constants, are called initial conditions. Those of you who are familiar with differential equations may recognize that the this language parallels what is used in differential equations. Difference (AKA recurrence) equations move forward discretely—that is, in a finite number of positive steps—while a differential equation moves continuously—that is, takes an infinite number of infinitesimal steps.

(2) A recurrence relationship of the form $F_k = aF_{k-1} + b$, where a and b are constants, is called a first-order difference equation. In order to write out the sequence, we need to know one initial condition. Equations of this type can be solved similarly to the method outlined in Example 12.5.1 by introducing the superfluous equation $1 = 0F_{k-1} + 1$ to obtain in matrix equation:

$$\begin{pmatrix} F_k \\ 1 \end{pmatrix} = \begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix} \begin{pmatrix} F_{k-1} \\ 1 \end{pmatrix} \Rightarrow \begin{pmatrix} F_k \\ 1 \end{pmatrix} = \begin{pmatrix} a & b \\ 0 & 1 \end{pmatrix}^k \begin{pmatrix} F_0 \\ 1 \end{pmatrix}$$

Example 12.5.2: Graph Theory. Consider the graph in Figure 12.5.1.

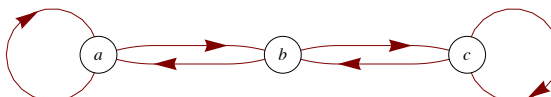


Figure 12.5.1

From the procedures outlined in Section 6.4, the adjacency matrix of this graph is

$$A = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix}$$

Recall that A^k is the adjacency matrix of the relation r^k , where r is the relation $\{(a, a), (a, b), (b, a), (b, c), (c, b), (c, c)\}$ of the above graph. Also recall that in computing A^k , we used Boolean arithmetic. What happens if we use "regular" arithmetic? For example,

$$A = \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}$$

How can we interpret this? We note that $A_{33} = 2$ and that there are two paths of length two from c (the third node) to c . Also, $A_{13} = 1$, and there is one path of length 2 from a to c . The reader should verify these claims from the graph in Figure 12.5.1.

Theorem 12.5.1. The entry $(A^k)_{ij}$ is the number of paths, or walks, of length k from node v_i , to node v_j .

How do we find A^k for possibly large values of k ? From the discussion at the beginning of this section, we know that $A^k = P D^k P^{-1}$ if A is diagonalizable. We leave to the reader to show that $\lambda = 1, 2$, and -1 are eigenvalues of A with eigenvectors

$$\begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \text{ and } \begin{pmatrix} 1 \\ -2 \\ 1 \end{pmatrix}$$

respectively, so that

$$A^k = P \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2^k & 0 \\ 0 & 0 & (-1)^k \end{pmatrix} P^{-1}$$

$$\text{where } P = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & -2 \\ -1 & 1 & 1 \end{pmatrix} \text{ and } P^{-1} = \begin{pmatrix} \frac{1}{2} & 0 & -\frac{1}{2} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{6} & -\frac{1}{3} & \frac{1}{6} \end{pmatrix}$$

See Exercise 5 of this section for the completion of this example.

Example 12.5.3: Matrix Calculus. Those who have studied calculus recall that the Maclaurin series is a useful way of expressing many common functions. For example,

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

Indeed, calculators and computers use these series for calculations. Given a polynomial $f(x)$, we defined the matrix-polynomial $f(A)$ for square matrices in Chapter 5. Hence, we are in a position to describe e^A for an $n \times n$ matrix A as a limit of polynomial. Formally, we write

$$e^A = I + A + \frac{A^2}{2!} + \frac{A^3}{3!} + \cdots = \sum_{k=0}^{\infty} \frac{A^k}{k!}$$

Again we encounter the need to compute high powers of a matrix. Let A be an $n \times n$ diagonalizable matrix. Then there exists an invertible $n \times n$ matrix P such that $P^{-1}AP = D$, a diagonal matrix, so that

$$\begin{aligned} e^A &= e^{PDP^{-1}} \\ &= \sum_{k=0}^{\infty} \frac{(PDP^{-1})^k}{k!} \\ &= P \left(\sum_{k=0}^{\infty} \frac{D^k}{k!} \right) P^{-1} \end{aligned}$$

The infinite sum in the middle of this final expression can be easily evaluated if D is diagonal. All entries of powers off the diagonal are zero and the i^{th} entry of the diagonal is

$$\left(\sum_{k=0}^{\infty} \frac{D^k}{k!} \right)_{ii} = \sum_{k=0}^{\infty} \frac{D_{ii}^k}{k!} = e^{D_{ii}}$$

For example, if $A = \begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix}$, the first matrix we diagonalized in Section 12.3, we found that $P = \begin{pmatrix} 1 & 1 \\ -1 & 2 \end{pmatrix}$ and $D = \begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix}$. Therefore,

$$\begin{aligned} e^A &= \begin{pmatrix} 1 & 1 \\ -1 & 2 \end{pmatrix} \begin{pmatrix} e & 0 \\ 0 & e^4 \end{pmatrix} \begin{pmatrix} \frac{2}{3} & -\frac{1}{3} \\ \frac{1}{3} & \frac{1}{3} \end{pmatrix} \\ &= \begin{pmatrix} \frac{2e}{3} + \frac{e^4}{3} & -\frac{e}{3} + \frac{e^4}{3} \\ -\frac{2e}{3} + \frac{2e^4}{3} & \frac{e}{3} + \frac{2e^4}{3} \end{pmatrix} \\ &\approx \begin{pmatrix} 20.0116 & 17.2933 \\ 34.5866 & 37.3049 \end{pmatrix} \end{aligned}$$

Comments on Example 12.5.3:

(1) Many of the ideas of calculus can be developed using matrices. For example, if

$$A(t) = \begin{pmatrix} t^3 & 3t^2 + 8t \\ e^t & 2 \end{pmatrix}$$

then

$$\frac{dA(t)}{dt} = \begin{pmatrix} 3t^2 & 6t + 8 \\ e^t & 0 \end{pmatrix}$$

(2) Many of the basic formulas in calculus are true in matrix calculus. For example,

$$\frac{d(A(t)+B(t))}{dt} = \frac{dA(t)}{dt} + \frac{dB(t)}{dt}$$

and if A is a constant matrix,

$$\frac{d e^{At}}{dt} = A e^{At}$$

(3) Matrix calculus can be used to solve systems of differential equations in a similar manner to the procedure used in ordinary differential equations.



Mathematica Note

Mathematica's matrix exponential function is **MatrixExp**.

$$\text{MatrixExp}\left[\begin{pmatrix} 2 & 1 \\ 2 & 3 \end{pmatrix}\right]$$

$$\begin{pmatrix} \frac{1}{3}(2e + e^4) & \frac{1}{3}(-e + e^4) \\ \frac{2}{3}(-e + e^4) & \frac{1}{3}(e + 2e^4) \end{pmatrix}$$



Sage Note

Sage's matrix exponential method is called **exp**.

```
A=Matrix(QQ,[[2,1],[2,3]]);
A.exp()
[ 2/3*e + 1/3*e^4 -1/3*e + 1/3*e^4]
[-2/3*e + 2/3*e^4  1/3*e + 2/3*e^4]
```

EXERCISES FOR SECTION 12.5

A Exercises

- (a) Write out all the details of Example 12.5.1 to show that the formula for F_k given in the text is correct.
(b) Use induction to prove the assertion made in Example 12.5.1 that

$$\begin{pmatrix} F_k \\ F_{k-1} \end{pmatrix} = A^{k-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

- (a) Do Example 8.3.8 of Chapter 8 using the method outlined in Example 12.5.1. Note that the terminology characteristic equation, characteristic polynomial, and so on, introduced in Chapter 8, comes from the language of matrix algebra.
(b) What is the significance of Algorithm 8.3.1, part c, with respect to this section?
- Solve $S(k) = 5S(k-1) + 4$, with $S(0) = 0$, using the method of this section.
- How many paths are there of length 6 between vertex 1 and vertex 3 in Figure 12.5.2? How many paths from vertex 2 to vertex 2 of length 6 are there? Hint: The characteristic polynomial of the adjacency matrix is λ^4 .

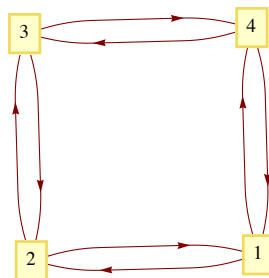


Figure 12.5.2

- Use the matrix A of Example 12.5.2 to:
 - Determine the number of paths of length 1 that exist from vertex a to each of the vertices in Example 12.5.2. Verify using the graph. Do the same for vertices b and c .
 - Verify all the details of Example 12.5.2.

(c) Use Example 12.5.2 to determine the number of paths of length 4 there are from each node in the graph of Figure 12.5.1 to every node in the graph. Verify your results using the graph.

6. Let $A = \begin{pmatrix} 2 & -1 \\ -1 & 2 \end{pmatrix}$

(a) Find e^A

(b) Recall that $\sin x = \sum_{k=0}^{\infty} \frac{(-1)^k x^{2k+1}}{(2k+1)!}$ and compute $\sin A$.

(d) Formulate a reasonable definition of the natural logarithm of a matrix and compute $\ln A$.

7. We noted in Chapter 5 that since matrix algebra is not commutative under multiplication, certain difficulties arise. Let $A = \begin{pmatrix} 1 & 1 \\ 0 & 0 \end{pmatrix}$ and $B = \begin{pmatrix} 0 & 0 \\ 0 & 2 \end{pmatrix}$.

(a) Compute e^A , e^B , and e^{A+B} . Compare $e^A e^B$, $e^B e^A$ and e^{A+B} .

(b) Show that if $\mathbf{0}$ is the 2×2 zero matrix, then $e^{\mathbf{0}} = I$.

(c) Prove that if A and B are two matrices that do commute, then $e^{A+B} = e^A e^B$, thereby proving that e^A and e^B commute.

(d) Prove that for any matrix A , $(e^A)^{-1} = e^{-A}$.

8. Another observation for adjacency matrices: For the matrix in Example 12.5.2, note that the sum of the elements in the row corresponding to the node a (that is, the first row) gives the outdegree of a . Similarly, the sum of the elements in any given column gives the indegree of the node corresponding to that column.

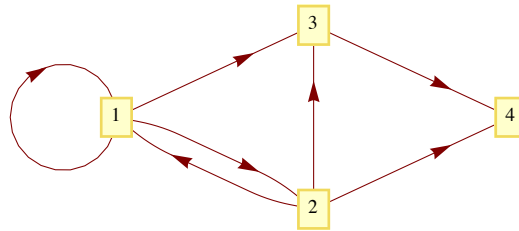


Figure 12.5.3

(a) Using the matrix A of Example 12.5.2, find the outdegree and the indegree of each node. Verify by the graph.

(b) Repeat part (a) for the directed graphs in Figure 12.5.3.

SUPPLEMENTARY EXERCISES FOR CHAPTER 12

Section 12.1

1. Find all solutions of the following systems:

$$\begin{array}{ll} \text{(a)} & \begin{array}{l} 2x_1 - 2x_2 + x_3 = 1 \\ x_2 - x_3 = 0 \\ x_1 + x_2 + x_3 = 3 \end{array} \\ \text{(b)} & \begin{array}{l} x_1 - x_3 = 0 \\ 2x_1 - 4x_2 = 1 \\ -x_1 + x_2 - x_3 = -1 \end{array} \end{array}$$

2. Find all solutions of

$$\begin{array}{l} x_1 - x_2 + 2x_3 = 1 \\ 3x_1 + x_3 = 2 \\ 2x_1 + x_2 - x_3 = 1 \end{array}$$

Section 12.2

3. Determine A^{-1} using the method of the text if

$$A = \begin{pmatrix} 1 & 2 & 1 \\ -2 & -3 & -1 \\ 1 & 4 & 4 \end{pmatrix}.$$

4. Find the inverse of the matrix

$$\begin{pmatrix} 0 & -4 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 4 \end{pmatrix}.$$

Section 12.3

5. In this exercise, write elements of \mathbb{R}^2 in column form. Let $\{x_1, x_2\}$ be a basis in \mathbb{R}^2 . Prove that $\{Ax_1, Ax_2\}$ is a basis for \mathbb{R}^2 if and only if A has an inverse.
6. Let $V = \{f : X \rightarrow \mathbb{R}\}$, where X is any nonempty set. Show that V is a vector space under the operations:

$$(f + g)(x) = f(x) + g(x) \text{ for } f, g \in V, \text{ and } x \in X$$

$$(cf)(x) = cf(x) \text{ for } f \in V, c \in \mathbb{R}, \text{ and } x \in X.$$

7. (a) Convince yourself that $M_{2 \times 3}(\mathbb{Z}_2)$ is a vector space over \mathbb{Z}_2 (i.e., allow only scalars from \mathbb{Z}_2 and use mod 2 arithmetic).
 (b) What is the vector $-\mathbf{X}$, for any $\mathbf{X} \in M_{2 \times 3}(\mathbb{Z}_2)$?
 (c) What is $|M_{2 \times 3}(\mathbb{Z}_2)|$?
8. (a) Define operations on \mathbb{R} so that \mathbb{R} is a vector space over \mathbb{R} .
 (b) What is a basis for the vector space part a? What is its dimension?

Section 12.4

9. Employ the diagonalization process to approximate the 100th power of A , where $A = \begin{bmatrix} 0.6 & 0.2 \\ 0.4 & 0.8 \end{bmatrix}$.

10. Let $B = \begin{pmatrix} 0 & -\frac{3}{5} & 0 \\ \frac{5}{3} & 0 & -\frac{5}{3} \\ 0 & 6 & -6 \end{pmatrix}$ and $C = \begin{pmatrix} 4 & 2 & 2 \\ 2 & 4 & 2 \\ 2 & 2 & 4 \end{pmatrix}$

- (a) Find all of the eigenvalues of B .

- (b) Given that 2 and 8 are the only eigenvalues of C , find invertible matrix P and diagonal matrix D such that $C = PDP^{-1}$.

11. Let $A = \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$ and $B = \begin{pmatrix} 4 & 0 & 0 \\ 0 & 4 & 1 \\ 0 & 0 & 2 \end{pmatrix}$

- (a) Find all of the eigenvalues of A .
- (b) Given that 4 and 2 are the only eigenvalues of B , find invertible matrix P and diagonal matrix D such that $B = PDP^{-1}$.
12. Find all eigenvalues and associated eigenvectors of the matrix A , and write A in the form $A = PDP^{-1}$.

$$A = \begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix}$$

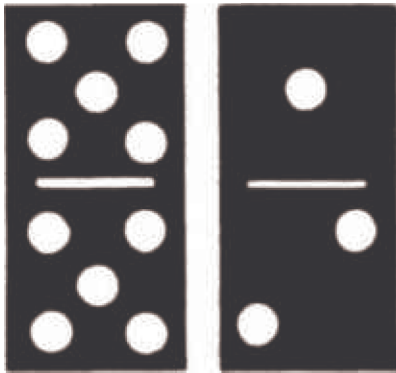
Section 12.5

13. For a multigraph we can define its matrix representation as follows: A_{ij} = the number of different edges e from vertex a_i to vertex a_j .
- (a) Draw the digraph that is described by the following matrix:

$$A = \begin{pmatrix} 2 & 1 & 0 \\ 1 & 0 & 3 \\ 1 & 1 & 0 \end{pmatrix}$$

- (b) Determine A^2 and interpret the result using Theorem 12.5.1.

Chapter 13



BOOLEAN ALGEBRA

GOALS

In this section we will develop an algebra that is particularly important to computer scientists, as it is the mathematical foundation of computer design, or switching theory. This algebra is called Boolean algebra after the mathematician George Boole (1815-64). The similarities of Boolean algebra and the algebra of sets and logic will be discussed, and we will discover special properties of finite Boolean algebras.



George Boole, 1815 - 1864

In order to achieve these goals, we will recall the basic ideas of posets introduced in Chapter 6 and develop the concept of a lattice, which has applications in finite-state machines.

The reader should view the development of the topics of this chapter as another example of an algebraic system. Hence, we expect to define first the elements in the system, next the operations on the elements, and then the common properties of the operations in the system.

13.1 Posets Revisited

From Chapter 6, Section 3, we recall the following definition:

Definition: *Poset.* A set L on which a partial ordering relation (reflexive, antisymmetric, and transitive) r is defined is called a partially ordered set, or poset, for short.

We recall a few examples of posets:

(1) $L = \mathbb{R}$ and r is the relation \leq .

(2) $L = \mathcal{P}(A)$ where $A = \{a, b\}$ and r is the relation \subseteq .

(3) $L = \{1, 2, 3, 6\}$ and r is the relation $|$ (divides). We remind the reader that the pair (a, b) as an element of the relation r can be expressed as $(a, b) \in r$, or $a r b$, depending on convenience and readability.

The posets we will concentrate on in this chapter will be those which have maxima and minima. These partial orderings resemble that of \leq on \mathbb{R} , so the symbol \leq is used to replace the symbol r in the definition of a partially ordered set. Hence, the definition of a poset becomes:

Definition: *Poset.* A set on which a partial ordering, \leq , is defined is called a partially ordered set, or, in brief, a poset. Here, \leq is a partial ordering on L if and only if for all $a, b, c \in L$:

(1) $a \leq a$ (reflexivity),

(2) $a \leq b$ and $b \leq a \Rightarrow a = b$ (antisymmetry), and

We now proceed to introduce maximum and minimum concepts. To do this, we will first define these concepts for two elements of the poset L , and then define the concepts over the whole poset L .

Definition: *Lower Bound, Upper Bound.* Let $a, b \in L$, a poset. Then $c \in L$ is a lower bound of a and b if $c \leq a$ and $c \leq b$. $d \in L$ is an upper bound of a and b if $a \leq d$ and $b \leq d$.

Definition: *Greatest Lower Bound.* Let L be a poset and \leq be the partial ordering on L . Let $a, b \in L$, then $g \in L$ is a greatest lower bound of a and b , denoted $\text{glb}(a, b)$, if and only if

- $g \leq a$,
- $g \leq b$, and
- if $g' \in L$ such that if $g' \leq a$ and $g' \leq b$, then $g' \leq g$.

The last condition says, in other words, that if g' is also a lower bound, then g is "greater" than g' , so g is a greatest lower bound.

The definition of a least upper bound is a mirror image of a greatest lower bound:

Definition: *Least Upper Bound.* Let L be a poset and \leq be the partial ordering on L . Let $a, b \in L$, then $\ell \in L$ is a least upper bound of a and b , denoted $\text{lub}(a, b)$, if and only if

- $a \leq \ell$,
- $b \leq \ell$, and
- if $\ell' \in L$ such that if $a \leq \ell'$ and $b \leq \ell'$, then $\ell \leq \ell'$.

Notice that the two definitions above refer to "...a greatest lower bound" and "a least upper bound." Any time you define an object like these you need to have an open mind as to whether more than one such object can exist. In fact, we now can prove that there can't be two greatest lower bounds or two least upper bounds.

Theorem 13.1.1. Let L be a poset and \leq be the partial ordering on L , and $a, b \in L$. If a greatest lower bound of a and b exists, then it is unique. The same is true of a least upper bound, if it exists.

Proof: Let g and g' be greatest lower bounds of a and b . We will prove that $g = g'$.

- (1) g a greatest lower bound of a and $b \Rightarrow g$ is a lower bound of a and b .
- (2) g' a greatest lower bound of a and b and g a lower bound of a and $b \Rightarrow g \leq g'$ by the definition of greatest lower bound.
- (3) g' a greatest lower bound of a and $b \Rightarrow g'$ is a lower bound of a and b .
- (4) g a greatest lower bound of a and b and g' a lower bound of a and $b \Rightarrow g' \leq g$ by the definition of greatest lower bound.
- (5) $g \leq g'$ and $g' \leq g \Rightarrow g = g'$ by the antisymmetry property of a partial ordering.

The proof of the second statement in the theorem is almost identical to the first and is left to the reader. ■

Definition: *Greatest Element, Least Element.* Let L be a poset. $M \in L$ is called the greatest (maximum) element of L if, for all $a \in L$, $a \leq M$. In addition, $m \in L$ is called the least (minimum) element of L if for all $a \in L$, $m \leq a$.

Note: The greatest and least elements, when they exist, are frequently denoted by 1 and 0 respectively.

Example 13.1.1. Let $L = \{1, 3, 5, 7, 15, 21, 35, 105\}$ and let \leq be the relation $|$ (divides) on L . Then L is a poset. To determine the lub of 3 and 7, we look for all $\ell \in L$, such that $3 | \ell$ and $7 | \ell$. Certainly, both $\ell = 21$ and $\ell = 105$ satisfy these conditions and no other element of L does. Next, since $21 | 105$, then $21 = \text{lub}(3, 7)$. Similarly, the $\text{lub}(3, 5) = 15$. The greatest element of L is 105 since $a | 105$ for all $a \in L$. To find the glb of 15 and 35, we first consider all elements g of L such that $g | 15$ and $g | 35$. Certainly, both $g = 5$ and $g = 1$ satisfy these conditions. But since $1 | 5$, then $\text{glb}(15, 35) = 5$. The least element of L is 1 since $1 | a$ for all $a \in L$.

Henceforth, for any positive integer n , D_n will denote the set of all positive integers which are divisors of n . For example, the set L of Example 13.1.1 is D_{105} .

Example 13.1.2. Consider the poset $\mathcal{P}(A)$, where $A = \{a, b, c\}$, with the relation \subseteq on $\mathcal{P}(A)$. The glb of the $\{a, b\}$ and $\{a, c\}$ is $g = \{a\}$. For any other element g' of M which is a subset of $\{a, b\}$ and $\{a, c\}$ (there is only one; what is it?), $g' \subseteq g$. The least element of $\mathcal{P}(A)$ is \emptyset and the greatest element of $\mathcal{P}(A)$ is $A = \{a, b, c\}$. The Hasse diagram of $\mathcal{P}(A)$ is shown in Figure 13.1.1.

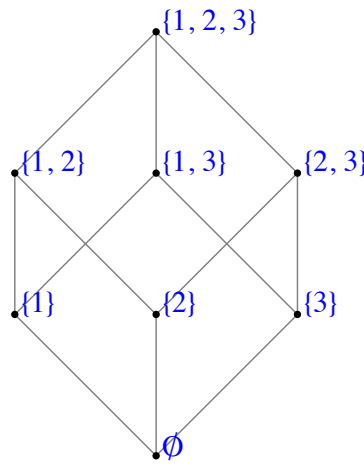


Figure 13.1.1
Example 13.1.2

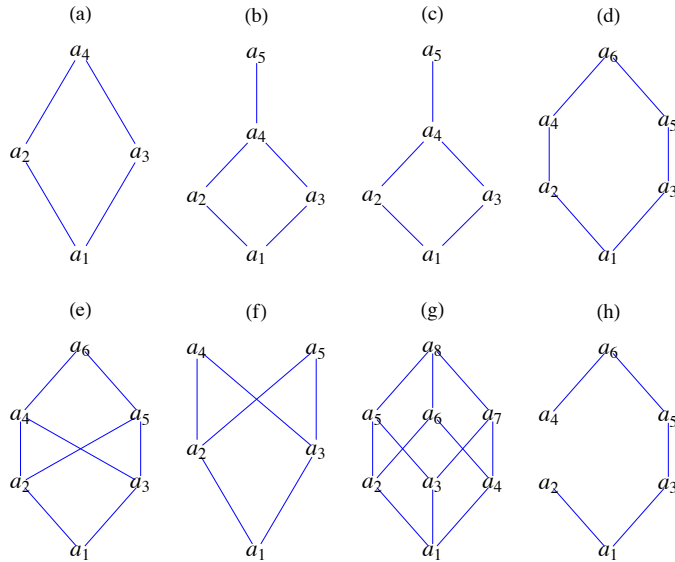
With a little practice, it is quite easy to find the least upper bounds and greatest lower bounds of all possible pairs in $\mathcal{P}(A)$ directly from the graph of the poset.

The previous examples and definitions indicate that the lub and glb are defined in terms of the partial ordering of the given poset. It is not yet clear whether all posets have the property such every pair of elements has both a lub and a glb . Indeed, this is not the case (see Exercise 3).

EXERCISES FOR SECTION 13.1

A Exercises

- Let $D_{30} = \{1, 2, 3, 5, 6, 10, 15, 30\}$ and let the relation $|$ be a partial ordering on D_{30} .
 - Find all lower bounds of 10 and 15.
 - Find the glb of 10 and 15.
 - Find all upper bounds of 10 and 15.
 - Determine the lub of 10 and 15.
 - Draw the Hasse diagram for D_{30} with $|$. Compare this Hasse diagram with that of Example 13.1.2. Note that the two diagrams are structurally the same.
- List the elements of the sets D_8 , D_{50} , and D_{1001} . For each set, draw the Hasse diagram for "divides."
- Figure 13.1.2 contains Hasse diagrams of posets.
 - Determine the lub and glb of all pairs of elements when they exist. Indicate those pairs that do not have a lub (or a glb).
 - Find the least and greatest elements when they exist.

Figure 13.1.2
Exercise 3

4. For the poset (\mathbb{N}, \leq) , what are $\text{glb}(a, b)$ and $\text{lub}(a, b)$? Are there least and/or greatest elements?
5. (a) Prove the second part of Theorem 13.1.1, the least upper bound of two elements in a poset is unique, if one exists.
(b) Prove that if a poset L has a least element, then that element is unique.
6. We naturally order the numbers in $A_m = \{1, 2, \dots, m\}$ with "less than or equal to," which is a partial ordering. We may order the elements of $A_m \times A_n$ by $(a, b) \leq (a', b') \Leftrightarrow a \leq a' \text{ and } b \leq b'$.
(a) Prove that this defines a partial ordering of $A_m \times A_n$.
(b) Draw the ordering diagrams for \leq on $A_2 \times A_2$, $A_2 \times A_3$, and $A_3 \times A_3$.
(c) What are $\text{glb}((a, b), (a', b'))$ and $\text{lub}((a, b), (a', b'))$?
(d) Are there least and/or greatest elements in $A_m \times A_n$?

13.2 Lattices

In this section, we restrict our discussion to *lattices*, those posets where every pair of elements has a *lub* and a *glb*. We first introduce some notation.

Definitions: Join, Meet. Let L be a poset under an ordering \leq . Let $a, b \in L$. We define:

$a \vee b$ (read "a join b") as the least upper bound of a and b , and

$a \wedge b$ (read "a meet b") as greatest lower bound of a and b .

Since the join and meet operations produce a unique result in all cases where they exist, by Theorem 13.1.1, we can consider them as binary operations on a set if they always exist. Thus the following definition:

Definition: Lattice. A lattice is a poset L (under \leq) in which every pair of elements has a *lub* and a *glb*. Since a lattice L is an algebraic system with binary operations \vee and \wedge , it is denoted by $[L; \vee, \wedge]$.

In Example 13.1.2, the operation table for the *lub* operation is easy, although admittedly tedious, to do. We can observe that every pair of elements in this poset has a least upper bound. In fact, $A \vee B = A \cup B$.

The reader is encouraged to write out the operation table for the *glb* operation and to note that every pair of elements in this poset also has a *glb*, so that $\mathcal{P}(A)$ together with these two operations is a lattice. We further observe that:

- (1) $[\mathcal{P}(A); \vee, \wedge]$ is a lattice (under \subseteq) for any set A , and
- (2) the join operation is the set operation of union and the meet operation is the operation intersection; that is, $\vee = \cup$ and $\wedge = \cap$.

It can be shown (see the exercises) that the commutative laws, associative laws, idempotent laws, and absorption laws are all true for any lattice. An example of this is clearly $[\mathcal{P}(A); \cup, \cap]$, since these laws hold in the algebra of sets. This lattice is also distributive in that join is distributive over meet and meet is distributive over join. This is not always the case for lattices in general however.

Definition: Distributive Lattice. Let $[L; \vee, \wedge]$ be a lattice (under \leq). $[L; \vee, \wedge]$ is called a distributive lattice if and only if the distributive laws hold; that is, for all $a, b, c \in L$, we have:

$$a \vee (b \wedge c) = (a \vee b) \wedge (a \vee c) \text{ and}$$

$$a \wedge (b \vee c) = (a \wedge b) \vee (a \wedge c).$$

Example 13.2.1. If A is any set, the lattice $[\mathcal{P}(A); \cup, \cap]$ is distributive.

Example 13.2.2. We now give an example of a lattice where the distributive laws do not hold. Let $L = \{1, 2, 3, 5, 30\}$. Then L is a poset under the relation divides. The operation tables for \vee and \wedge on L are:

\vee	1	2	3	5	30
1	1	2	3	5	30
2	2	2	30	30	30
3	3	30	3	30	30
5	5	30	30	5	30
30	30	30	30	30	30

\wedge	1	2	3	5	30
1	1	1	1	1	1
2	1	2	1	1	2
3	1	1	3	1	3
5	1	1	1	5	5
30	1	2	3	5	30

Since every pair of elements in L has both a join and a meet, $[L; \vee, \wedge]$ is a lattice (under divides). Is this lattice distributive? We note that:

$$2 \vee (5 \wedge 3) = 2 \vee 1 = 2 \text{ and}$$

$$(2 \vee 5) \wedge (2 \vee 3) = 30 \wedge 30 = 30,$$

so that $a \vee (b \wedge c) \neq (a \vee b) \wedge (a \vee c)$ for some values of $a, b, c \in L$. Hence L is not a distributive lattice.

It can be shown that a lattice is nondistributive if and only if it contains a sublattice isomorphic to one of the lattices in Figure 13.2.1.

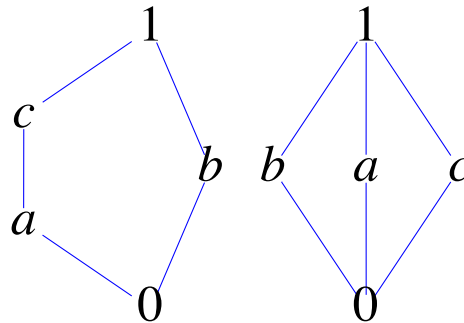


Figure 13.2.1
Nondistributive lattices

It is interesting to note that for the relation "divides" on \mathbb{P} , if $a, b \in \mathbb{P}$ we have:

$a \vee b = \text{lcm}(a, b)$, the least common multiple of a and b ; that is, the smallest integer (in \mathbb{P}) that is divisible by both a and b ;

$a \wedge b = \text{gcd}(a, b)$, the greatest common divisor of a and b ; that is, the largest integer that divides both a and b .

EXERCISES FOR SECTION 13.2

A Exercises

- Let L be the set of all propositions generated by p and q . What are the meet and join operations in this lattice. What are the maximum and minimum elements?
- Which of the posets in Exercise 3 of Section 13.1 are lattices? Which of the lattices are distributive?

B Exercises

- (a) State the commutative laws, associative laws, idempotent laws, and absorption laws for lattices.

- (b) Prove these laws.
4. Let $[L; \vee, \wedge]$ be a lattice based on a partial ordering \leq . Prove that if $a, b, c \in L$,
- (a) $a \vee b \geq a$.
- (b) $a \wedge b \leq a$.
- (c) $a \geq b$ and $a \geq c \Rightarrow a \geq b \vee c$.

13.3 Boolean Algebras

In order to define a Boolean algebra, we need the additional concept of complementation.

Definition: Complemented Lattice. Let $[L; \vee, \wedge]$ be a lattice that contains a least element, 0, and a greatest element, 1. $[L; \vee, \wedge]$ is called a *complemented lattice* if and only if for every element $a \in L$, there exists an element \bar{a} in L such that $a \wedge \bar{a} = 0$ and $a \vee \bar{a} = 1$. Such an element \bar{a} is called a *complement* of the element a .

Example 13.3.1. Let $L = \mathcal{P}(A)$, where $A = \{a, b, c\}$. Then $[L; \cup, \cap]$ is a bounded lattice with $0 = \emptyset$ and $1 = A$. Then, to find if it exists, the complement, \bar{B} , of, say $B = \{a, b\} \in L$, we want \bar{B} such that

$$\{a, b\} \cap \bar{B} = \emptyset \text{ and } \{a, b\} \cup \bar{B} = A.$$

Here, $\bar{B} = \{c\}$, and since it can be shown that each element of L has a complement (see Exercise 1), $[L; \cup, \cap]$ is a complemented lattice. Note that if A is any set and $L = \mathcal{P}(A)$, then $[L; \cup, \cap]$ is a complemented lattice where the complement of $B \in L$ is $\bar{B} = B^c = A - B$.

In Example 13.3.1, we observe that the complement of each element of L is unique. Is this always the case? The answer is no. Consider the following.

Example 13.3.2. Let $L = \{1, 2, 3, 5, 30\}$ and consider the lattice $[L; \vee, \wedge]$ (under "divides"). The least element of L is 1 and the greatest element is 30. Let us compute the complement of the element $a = 2$. We want to determine \bar{a} such that $2 \wedge \bar{a} = 1$ and $2 \vee \bar{a} = 30$. Certainly, $\bar{a} = 3$ works, but so does $\bar{a} = 5$, so the complement of $a = 2$ in this lattice is not unique. However, $[L; \vee, \wedge]$ is still a complemented lattice since each element does have at least one complement.

The following theorem gives us an insight into when uniqueness of complements occurs.

Theorem 13.3.1. If $[L; \vee, \wedge]$ is a complemented and distributive lattice, then the complement \bar{a} of any element $a \in L$ is unique.

Proof: Let $a \in L$ and assume to the contrary that a has two complements, namely a_1 and a_2 . Then by definition of complement,

$$a \wedge a_1 = 0 \text{ and } a \vee a_1 = 1,$$

Also,

$$a \wedge a_2 = 0 \text{ and } a \vee a_2 = 1.$$

So that

$$\begin{aligned} a_1 &= a_1 \wedge 1 = a_1 \wedge (a \vee a_2) \\ &= (a_1 \wedge a) \vee (a_1 \wedge a_2) \\ &= 0 \vee (a_1 \wedge a_2) \\ &= a_1 \wedge a_2. \end{aligned}$$

On the other hand,

$$\begin{aligned} a_2 &= a_2 \wedge 1 = a_2 \wedge (a \vee a_1) \\ &= (a_2 \wedge a) \vee (a_2 \wedge a_1) \\ &= 0 \vee (a_2 \wedge a_1) \\ &= a_2 \wedge a_1. \end{aligned}$$

Hence $a_1 = a_2$, which contradicts the assumption that a has two different complements, a_1 and a_2 . ■

Definition: Boolean Algebra. A Boolean algebra is a lattice that contains a least element and a greatest element and that is both complemented and distributive.

Since the complement of each element in a Boolean algebra is unique (by Theorem 13.3.1), complementation is a valid unary operation over the set under discussion, and we will list it together with the other two operations to emphasize that we are discussing a set together with three operations. Also, to help emphasize the distinction between lattices and lattices that are Boolean algebras, we will use the letter B as the generic symbol for the set of a Boolean algebra; that is, $[B; -, \vee, \wedge]$ will stand for a general Boolean algebra.

Example 13.3.3. Let A be any set, and let $B = \mathcal{P}(A)$. Then $[B; c, \cup, \cap]$ is a Boolean algebra. Here, c stands for the complement of an element of B with respect to A , $A - B$.

This is a key example for us since all finite Boolean algebras and many infinite Boolean algebras look like this example for some A . In fact, a glance at the basic Boolean algebra laws in Table 13.3.1, in comparison with the set laws of Chapter 4 and the basic laws of logic of Chapter 3,

indicates that all three systems behave the same; that is, they are isomorphic.

The "pairing" of the above laws reminds us of the principle of duality, which we state for a Boolean algebra.

Definition: Principle of Duality for Boolean Algebras. Let $[B; -, \vee, \wedge]$ be a Boolean algebra (under \leq), and let S be a true statement for $[B; -, \vee, \wedge]$. If S^* is obtained from S by replacing \leq by \geq (this is equivalent to turning the graph upside down), \vee by \wedge , \wedge by \vee , 0 by 1 , and 1 by 0 , then S^* is also a true statement.

TABLE 13.3.1
Basic Boolean Algebra Laws

Commutative Laws	
1. $a \vee b = b \vee a$	1.' $a \wedge b = b \wedge a$
Associative Laws	
2. $a \vee (b \vee c) = (a \vee b) \vee c$	2.' $a \wedge (b \wedge c) = (a \wedge b) \wedge c$
Distributive Laws	
3. $a \wedge (b \vee c) = (a \wedge b) \vee (a \wedge c)$	3.' $a \vee (b \wedge c) = (a \vee b) \wedge (a \vee c)$
Identity Laws	
4. $a \vee 0 = 0 \vee a = a$	4.' $a \wedge 1 = 1 \wedge a = a$
Complement Laws	
5. $a \vee \bar{a} = 1$	5.' $a \wedge \bar{a} = 0$
Idempotent Laws	
6. $a \vee a = a$	6.' $a \wedge a = a$
Null Laws	
7. $a \vee 1 = 1$	7.' $a \wedge 0 = 0$
Absorption Laws	
8. $a \vee (a \wedge b) = a$	8.' $a \wedge (a \vee b) = a$
DeMorgan's Laws	
9. $\overline{a \vee b} = \bar{a} \wedge \bar{b}$	9.' $\overline{a \wedge b} = \bar{a} \vee \bar{b}$
Involution Law	
10. $\overline{\bar{a}} = a$	

Example 13.3.4. The laws 1' through 9' are the duals of the Laws 1 through 9 respectively. Law 10 is its own dual.

We close this section with some comments on notation. The notation for operations in a Boolean algebra is derived from the algebra of logic. However, other notations are used. These are summarized in the following chart;

Notation used in this text (Mathematics notation)	Set Notation	Logic Design (CS/EE notation)	Read as
\vee	\cup	\oplus	join
\wedge	\cap	\otimes	meet
$-$	c	$-$	complement
\leq	\subseteq	\leq	underlying partial ordering

Mathematicians most frequently use the notation of the text, and, on occasion, use set notation for Boolean algebras. Thinking in terms of sets may be easier for some people. Computer designers traditionally use the arithmetic and notation. In this latter notation, DeMorgan's Laws become:

$$(9) \quad \overline{a \oplus b} = \bar{a} \otimes \bar{b}$$

and

$$(9') \quad \overline{a \otimes b} = \bar{a} \oplus \bar{b}.$$

EXERCISES FOR SECTION 13.3

A Exercises

- Determine the complement of each element $B \in L$ in Example 13.3.1. Is this lattice a Boolean algebra? Why?
- Determine the complement of each element of D_6 in $[D_6; \vee, \wedge]$.
 - Repeat part a using the lattice in Example 13.2.2.
 - Repeat part a using the lattice in Exercise 1 of Section 13.1.
 - Are the lattices in parts a, b, and c Boolean algebras? Why?
- Determine which of the lattices of Exercise 3 of Section 13.1 are Boolean algebras.
- Let $A = \{a, b\}$ and $B = \mathcal{P}(A)$.
 - Prove that $[B; c, \cup, \cap]$ is a Boolean algebra.
 - Write out the operation tables for the Boolean algebra.
- It can be shown that the following statement, S , holds for any Boolean algebra $[B; -, \vee, \wedge] : (a \wedge b) = a$ if $a \leq b$.
 - Write the dual, S^* , of the statement S .
 - Write the statement S and its dual, S^* , in the language of sets.
 - Are the statements in part b true for all sets?
 - Write the statement S and its dual, S^* , in the language of logic.
 - Are the statements in part d true for all propositions?
- State the dual of:
 - $a \vee (b \wedge a) = a$.
 - $a \vee ((\bar{b} \vee a) \wedge b) = 1$.
 - $(a \wedge \bar{b}) \wedge b = a \vee b$.

B Exercises

- Formulate a definition for isomorphic Boolean algebras.

13.4 Atoms of a Boolean Algebra

In this section we will look more closely at previous claims that every finite Boolean algebra is isomorphic to an algebra of sets. We will show that every finite Boolean algebra has 2^n elements for some n with precisely n generators, called *atoms*.

Consider the Boolean algebra $[B; -, \vee, \wedge]$, whose graph is:

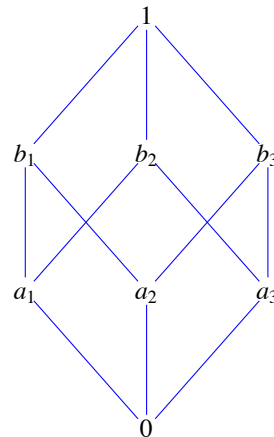


Figure 13.4.1
Illustration of the atom concept

We note that $1 = a_1 \vee a_2 \vee a_3$, $b_1 = a_1 \vee a_2$, $b_2 = a_1 \vee a_3$, and $b_3 = a_2 \vee a_3$; that is, each of the elements above level one can be described completely and uniquely in terms of the elements on level one. The a_i 's have uniquely generated the nonzero elements of B much like a basis in linear algebra generates the elements in a vector space. We also note that the a_i 's are the immediate successors of the minimum element, 0. In any Boolean algebra, the immediate successors of the minimum element are called *atoms*. Let A be any nonempty set. In the Boolean algebra $[\mathcal{P}(A); \subseteq, \cup, \cap]$ (over \subseteq), the singleton sets are the generators, or atoms, of the algebraic structure since each element $\mathcal{P}(A)$ can be described completely and uniquely as the join or union of singleton sets.

Definition: Atom. A nonzero element a in a Boolean algebra $[B; -, \vee, \wedge]$ is called an *atom* if for every $x \in B$, $x \wedge a = a$ or $x \wedge a = 0$.

The condition that $x \wedge a = a$ tells us that x is a successor of a ; that is, $a \leq x$, as depicted in Figure 13.4.2a.

The condition $x \wedge a = 0$ is true only when x and a are "not connected." This occurs when x is another atom or if x is a successor of atoms different from a , as depicted in Figure 13.4.2b.

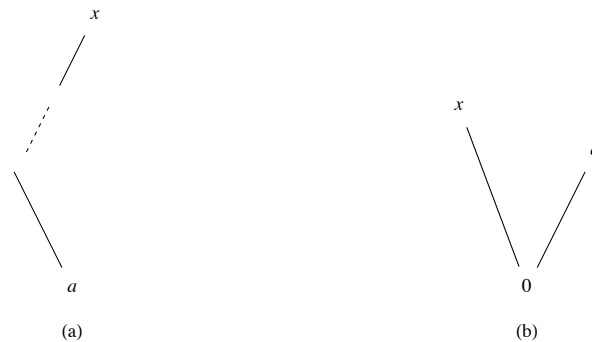


Figure 13.4.2

Example 13.4.1. The set of atoms of the Boolean algebra $[D_{30}; -, \vee, \wedge]$ is $M = \{2, 3, 5\}$. To see that $a = 2$ is an atom, let x be any nonzero element of D_{30} and note that one of the two conditions $x \wedge 2 = 2$ or $x \wedge 2 = 1$ holds. Of course, to apply the definition to this Boolean algebra, we must remind ourselves that in this case the 0-element is 1, the operation \wedge is *gcd*, and the poset relation \leq is "divides." So if $x = 10$, we have $10 \wedge 2 = 2$ (or $2 \mid 10$), so Condition 1 holds. If $x = 15$, the first condition is not true. (Why?) However, Condition 2, $15 \wedge 2 = 1$, is true. The reader is encouraged to show that each of the elements 2, 3, and 5 satisfy the definition (see Exercise 13.4.1). Next, if we compute the join (*lcm* in this case) of all possible combinations of the atoms 2, 3, and 5, we will generate all nonzero elements of D_{30} . For example, $2 \vee 3 \vee 5 = 30$ and $2 \vee 5 = 10$. We state this concept formally in the following theorem, which we give without proof.

Theorem 13.4.1. Let $[B; -, \vee, \wedge]$ be any finite Boolean algebra. Let $A = \{a_1, a_2, \dots, a_n\}$ be the set of all n atoms of $[B; -, \vee, \wedge]$. Then every nonzero element in B can be expressed uniquely as the join of a subset of A .

We now ask ourselves if we can be more definitive about the structure of different Boolean algebras of a given order. Certainly, the Boolean algebras $[D_{30}; -, \vee, \wedge]$ and $[\mathcal{P}(A); \subseteq, \cup, \cap]$ have the same graph (that of Figure 13.4.1), the same number of atoms, and, in all respects, look the same except for the names of the elements and the operations. In fact, when we apply corresponding operations to corresponding elements, we obtain corresponding results. We know from Chapter 11 that this means that the two structures are isomorphic as Boolean algebras. Furthermore, the graphs of these examples are exactly the same as that of Figure 13.4.1, which is an arbitrary Boolean algebra of order $8 = 2^3$.

In these examples of a Boolean algebra of order 8, we note that each had 3 atoms and $2^3 = 8$ number of elements, and all were isomorphic to $[\mathcal{P}(A); \subseteq, \cup, \cap]$, where $A = \{a, b, c\}$. This leads us to the following questions:

- (1) Are there any other different (nonisomorphic) Boolean algebras of order 8?
- (2) What is the relationship, if any, between finite Boolean algebras and their atoms?
- (3) How many different (nonisomorphic) Boolean algebras are there of order 2? Order 3? Order 4? And so on.

The answers to these questions are given in the following theorem and corollaries. We include the proofs of the corollaries since they are instructive.

Theorem 13.4.2. *Let $[B; -, \vee, \wedge]$ be any finite Boolean algebra, and let A be the set of all atoms in this Boolean algebra. Then $[B; -, \vee, \wedge]$ is isomorphic to $[\mathcal{P}(A); \subseteq, \cup, \cap]$.*

Corollary 13.4.1. *Every finite Boolean algebra $[B; -, \vee, \wedge]$ has 2^n elements for some positive integer n .*

Proof: Let A be the set of all atoms of B and let $|A| = n$. Then there are exactly 2^n elements (subsets) in $\mathcal{P}(A)$, and by Theorem 13.4.2, $[B; -, \vee, \wedge]$ is isomorphic to $[\mathcal{P}(A); \subseteq, \cup, \cap]$. ■

Corollary 13.4.2. All Boolean algebras of order 2^n are isomorphic to each other. (The graph of the Boolean algebra of order 2^n is the n -cube).

Proof: By Theorem 13.4.2, every Boolean algebra of order 2^n is isomorphic to $[\mathcal{P}(A); \subseteq, \cup, \cap]$ when $|A| = n$. Hence, they are all isomorphic to one another. ■

The above theorem and corollaries tell us that we can only have finite Boolean algebras of orders $2^1, 2^2, 2^3, \dots, 2^n$, and that all finite Boolean algebras of any given order are isomorphic. These are powerful tools in determining the structure of finite Boolean algebras. In the next section, we will try to find the easiest way of describing a Boolean algebra of any given order.

EXERCISES FOR SECTION 13.4

A Exercises

1. (a) Show that $a = 2$ is an atom of the Boolean algebra $[D_{30}; -, \vee, \wedge]$.
 (b) Repeat part a for the elements 3 and 5 of D_{30} .
 (c) Verify Theorem 13.4.1 for the Boolean algebra $[D_{30}; -, \vee, \wedge]$.
2. Let $A = \{a, b, c\}$.
 (a) Rewrite the definition of atom for $[\mathcal{P}(A); \subseteq, \cup, \cap]$. What does $a \leq x$ mean in this example?
 (b) Find all atoms of $[\mathcal{P}(A); \subseteq, \cup, \cap]$.
 (c) Verify Theorem 13.4.1 for $[\mathcal{P}(A); \subseteq, \cup, \cap]$.
3. Verify Theorem 13.4.2 and its corollaries for the Boolean algebras in Exercises 1 and 2 of this section.
4. Give a description of all Boolean algebras of order 16. (*Hint:* Use Theorem 13.4.2.) Note that the graph of this Boolean algebra is given in Figure 9.4.5.
5. Corollary 13.4.1 states that there do not exist Boolean algebras of orders 3, 5, 6, 7, 9, etc. (orders different from 2^n). Prove that we cannot have a Boolean algebra of order 3. (*Hint:* Assume that $[B; -, \vee, \wedge]$ is a Boolean algebra of order 3 where $B = \{0, x, 1\}$ and show that this cannot happen by investigating the possibilities for its operation tables.)
6. (a) There are many different, yet isomorphic, Boolean algebras with two elements. Describe one such Boolean algebra that is derived from a power set, $\mathcal{P}(A)$, under \subseteq . Describe a second that is described from D_n , for some $n \in \mathbb{P}$, under "divides."
 (b) Since the elements of a two-element Boolean algebra must be the greatest and least elements, 1 and 0, the tables for the operations on $\{0, 1\}$ are determined by the Boolean algebra laws. Write out the operation tables for $[\{0, 1\}; -, \vee, \wedge]$.

B Exercises

7. Find a Boolean algebra with a countably infinite number of elements.
8. Prove that the direct product of two Boolean algebras is a Boolean algebra. (*Hint:* "Copy" the corresponding proof for groups in Section 11.6.)

13.5 Finite Boolean Algebras as n-tuples of 0's and 1's

From the previous section we know that all finite Boolean algebras are of order 2^n , where n is the number of atoms in the algebra. We can therefore completely describe every finite Boolean algebra by the algebra of power sets. Is there a more convenient, or at least an alternate way, of defining finite Boolean algebras? In Chapter 11 we found that we could produce new groups by taking Cartesian products of previously known groups. We imitate this process for Boolean algebras.

The simplest nontrivial Boolean algebra is the Boolean algebra on the set $B_2 = \{0, 1\}$. The ordering on B_2 is the natural one, $0 \leq 0, 0 \leq 1, 1 \leq 1$. If we treat 0 and 1 as the truth values "false" and "true," respectively, we see that the Boolean operations \vee (join) and \wedge (meet) are nothing more than the logical connectives \vee (or) and \wedge (and). The Boolean operation, $-$, (complementation) is the logical \neg (negation). In fact, this is why the symbols $-$, \vee , and \wedge were chosen as the names of the Boolean operations. The operation tables for $[B_2; -, \vee, \wedge]$ are simply those of "or," "and," and "not," which we repeat here:

\vee	0	1
0	0	1
1	1	1

\wedge	0	1
0	0	0
1	0	1

u	\bar{u}
0	1
1	0

By Theorem 13.4.2 and its corollaries, all Boolean algebras of order 2 are isomorphic to this one.

We know that if we form $B_2 \times B_2 = B_2^2$ we obtain the set $\{(0, 0), (0, 1), (1, 0), (1, 1)\}$, a set of order 4. We define operations on B_2^2 the natural way, namely, componentwise, so that $(0, 1) \vee (1, 1) = (0 \vee 1, 1 \vee 1) = (1, 1)$, $(0, 1) \wedge (1, 1) = (0 \wedge 1, 1 \wedge 1) = (0, 1)$ and $(\overline{0, 1}) = (\bar{0}, \bar{1}) = (1, 0)$. We claim that B_2^2 is a Boolean algebra under the componentwise operations. Hence, $[B_2^2; -, \vee, \wedge]$ is a Boolean algebra of order 4. Since all Boolean algebras of order 4 are isomorphic to each other, we have found a simple way of describing all Boolean algebras of order 4.

It is quite clear that we can describe any Boolean algebra of order 8 by considering $B_2 \times B_2 \times B_2 = B_2^3$ and, in general, any Boolean algebra of order 2^n — that is, all finite Boolean algebras — by $B_2^n = B_2 \times B_2 \times \cdots \times B_2$ (n factors).

EXERCISES FOR SECTION 13.5

A Exercises

- (a) Write out the operation tables for $[B_2^2; -, \vee, \wedge]$.
 (b) Draw the Hasse diagram for $[B_2^2; -, \vee, \wedge]$ and compare your results with Figure 9.4.6.
 (c) Find the atoms of this Boolean algebra.
- (a) Write out the operation table for $[B_2^3; -, \vee, \wedge]$.
 (b) Draw the Hasse diagram for $[B_2^3; -, \vee, \wedge]$ and compare the results with Figure 9.4.6.
- (a) List all atoms of B_2^4 .
 (b) Describe the atoms of B_2^n $n \geq 1$.

B Exercise

- Theorem 13.4.2 tells us we can think of any finite Boolean algebra in terms of sets. In Chapter 4, Section 3, we defined the terms *minset* and *minset normal form*. Rephrase these definitions in the language of Boolean algebra. The generalization of minsets are called *minterms*.

13.6 Boolean Expressions

In this section, we will use our background from the previous sections and set theory to develop a procedure for simplifying Boolean expressions. This procedure has considerable application to the simplification of circuits in switching theory or logical design.

Definition: *Boolean Expression.* Let $[B; -, \vee, \wedge]$ be any Boolean algebra. Let x_1, x_2, \dots, x_k be variables in B ; that is, variables that can assume values from B . A Boolean expression generated by x_1, x_2, \dots, x_k is any valid combination of the x_i and the elements of B with the operations of meet, join, and complementation.

This definition, as expected, is the analog of the definition of a proposition generated by a set of propositions, presented in Section 3.2.

Each Boolean expression generated by k variables, $e(x_1, \dots, x_k)$, defines a function $f: B^k \rightarrow B$ where $f(a_1, \dots, a_k) = e(a_1, \dots, a_k)$. If B is a finite Boolean algebra, then there are a finite number of functions from B^k into B . Those functions that are defined in terms of Boolean expressions are called *Boolean functions*. As we will see, there is an infinite number of Boolean expressions that define each Boolean function. Naturally, the "shortest" of these expressions will be preferred. Since electronic circuits can be described as Boolean functions with $B = B_2$, this economization is quite useful.

Example 13.6.1. Consider any Boolean algebra $[B; -, \vee, \wedge]$ of order 2. How many functions $f : B^2 \rightarrow B$ are there? First, all Boolean algebras of order 2 are isomorphic to $[B_2; -, \vee, \wedge]$ so we want to determine the number of functions $f : B_2^2 \rightarrow B_2$. If we consider a Boolean function of two variables, x_1 and x_2 , we note that each variable has two possible values 0 and 1, so there are 2^2 ways of assigning these two values to the $k = 2$ variables. Hence, the table below has $2^2 = 4$ rows. So far we have a table such as that labeled 13.6.1.

x_1	x_2	$f(x_1, x_2)$
0	0	?
0	1	?
1	0	?
1	1	?

Table 13.6.1

General Form Of Boolean Function $f(x_1, x_2)$ of Example 13.6.1

How many possible different function values $f(x_1, x_2)$ can there be? To list a few: $f_1(x_1, x_2) = x_1$, $f_2(x_1, x_2) = x_2$, $f_3(x_1, x_2) = x_1 \vee x_2$, $f_4(x_1, x_2) = (x_1 \wedge \bar{x}_2) \vee x_2$, $f_5(x_1, x_2) = x_1 \wedge x_2 \vee \bar{x}_2$, etc. Each of these will give a table like that of Table 13.6.1. The tables for f_1 , and f_3 appear in Table 13.6.2.

x_1	x_2	$f_1(x_1, x_2)$	x_1	x_2	$f_3(x_1, x_2)$
0	0	0	0	0	0
0	1	0	0	1	1
1	0	1	1	0	1
1	1	1	1	1	1

Table 13.6.2

Boolean Functions f_1 and f_3 of Example 13.6.1

Two functions are different if and only if their tables (values) are different for at least one row. Of course by using the basic laws of Boolean algebra we can see that $f_3 = f_4$. Why? So if we simply list by brute force all "combinations" of x_1 and x_2 we will obtain unnecessary duplication. However, we note that for any combination of the variables x_1 , and x_2 there are only two possible values for $f(x_1, x_2)$, namely 0 or 1. Thus, we could write $2^4 = 16$ different functions on 2 variables.

Now let's count the number of different Boolean functions in a more general setting. We will consider two cases: first, when $B = B_2$, and second, when B is any finite Boolean algebra with 2^n elements.

Let $B = B_2$. Each function $f : B^k \rightarrow B$ is defined in terms of a table having 2^k rows. Therefore, since there are two possible images for each element of B^k , there are 2 raised to the 2^k , or 2^{2^k} different functions. *We claim that every one of these functions is a Boolean function.*

Now suppose that $|B| = 2^n > 2$. A function from B^k into B can still be defined in terms of a table. There are $|B|^k$ rows to each table and $|B|$ possible images for each row. Therefore, there are 2^n raised to the power 2^k different functions. If $n > 1$, then not every one of these functions is a Boolean function. Notice that in counting the numbers of functions we are applying the result of Exercise 5 of Section 7.1.

Since all Boolean algebras are isomorphic to a Boolean algebra of sets, the analogues of statements in sets are useful in Boolean algebras.

Definition: *Minterm.* A Boolean expression generated by x_1, x_2, \dots, x_k that has the form

$$\bigwedge_{i=1}^k y_i,$$

where each y_i may be either x_i or \bar{x}_i is called a minterm generated by x_1, x_2, \dots, x_k .

By a direct application of the Product Rule we see that there are 2^k different minterms generated by x_1, \dots, x_k .

Definition: *Minterm Normal Form.* A Boolean expression generated by x_1, \dots, x_k is in minterm normal form if it is the join of expressions of the form $a \wedge m$, where $a \in B$ and m is a minterm generated by x_1, \dots, x_k . That is, it is of the form

$$\bigvee_{j=1}^p (a_j \wedge m_j),$$

where $p = 2^k$ and m_1, m_2, \dots, m_p are the minterms generated by x_1, \dots, x_k .

If $B = B_2$, then each a_j in a minterm normal form is either 0 or 1. Therefore, $a_j \wedge m_j$ is either 0 or m_j .

Theorem 13.6.1. Let $e(x_1, \dots, x_k)$ be a Boolean expression over B . There exists a unique minterm normal form $M(x_1, \dots, x_k)$ that is equivalent to $e(x_1, \dots, x_k)$ in the sense that e and M define the same function from B^k into B .

The uniqueness in this theorem does not include the possible ordering of the minterms in M (commonly referred to as "uniqueness up to the order of minterms"). The proof of this theorem would be quite lengthy, and not very instructive, so we will leave it to the interested reader to attempt. The implications of the theorem are very interesting, however.

If $|B| = 2^n$, then there are 2^n raised to the 2^k different minterm normal forms. Since each different minterm normal form defines a different function, there are a like number of Boolean functions from B^k into B . If $B = B_2$, there are as many Boolean functions (2 raised to the 2^k) as there are functions from B^k into B , since there are 2 raised to the 2^n functions from B^k into B . The significance of this result is that any desired

function can be obtained using electronic circuits having 0 or 1 (off or on, positive or negative) values, but more complex, multivalued circuits would not have this flexibility.

We will close this section by examining minterm normal forms for expressions over B_2 , since they are a starting point for circuit economization.

Example 13.6.2. Consider the Boolean expression $f(x_1, x_2) = x_1 \vee \bar{x}_2$. One method of determining the minterm normal form of f is to think in terms of sets. Consider the diagram with the usual translation of notation in Figure 13.6.1. Then $f(x_1, x_2) = (\bar{x}_1 \wedge \bar{x}_2) \vee (x_1 \wedge \bar{x}_2) \vee (x_1 \wedge x_2)$.

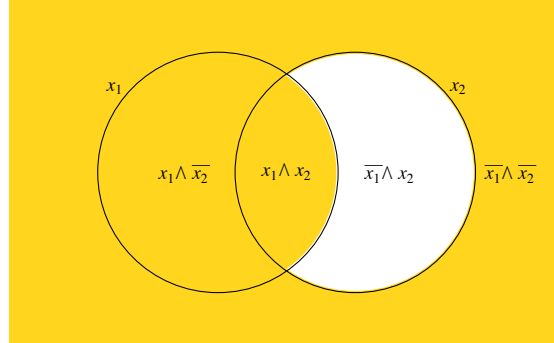


Figure 13.6.1

Example 13.6.3. Consider the function $f : B_2^3 \rightarrow B_2$ defined by Table 13.6.3. The minterm normal form for f can be obtained by taking the join of minterms that correspond to rows that have an image value of 1. If $f(a_1, a_2, a_3) = 1$, then include the minterm $y_1 \wedge y_2 \wedge y_3$ where

$$y_j = \begin{cases} x_j & \text{if } a_j = 1 \\ \bar{x}_j & \text{if } a_j = 0 \end{cases}$$

TABLE 13.6.3**Boolean Function of $f(a_1, a_2, a_3)$ Of Example 13.6.3**

a_1	a_2	a_3	$f(a_1, a_2, a_3)$
0	0	0	1
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	0
1	0	1	0
1	1	0	1
1	0	1	0

Therefore,

$$f(x_1, x_2, x_3) = (\bar{x}_1 \wedge \bar{x}_2 \wedge \bar{x}_3) \vee (\bar{x}_1 \wedge x_2 \wedge x_3) \vee (x_1 \wedge x_2 \wedge \bar{x}_3).$$

The minterm normal form is a first step in obtaining an economical way of expressing a given Boolean function. For functions of more than three variables, the above set theory approach tends to be awkward. Other procedures are used to write the normal form. The most convenient is the Karnaugh map, a discussion of which can be found in any logical design/switching theory text (see, for example, Hill and Peterson).

EXERCISES FOR SECTION 13.6

A Exercises

1. (a) Write the 16 possible functions of Example 13.6.1. (*Hint:* Find all possible joins of minterms generated by x_1 and x_2 .)
- (b) Write out the tables of several of the above Boolean functions to show that they are indeed different.
- (c) Determine the minterm normal form of

$$f_1(x_1, x_2) = x_1 \vee x_2,$$

$$f_2(x_1, x_2) = \bar{x}_1 \vee \bar{x}_2$$

$$f_3(x_1, x_2) = 0, f_4(x_1, x_2) = 1.$$

2. Consider the Boolean expression $f(x_1, x_2, x_3) = (\bar{x}_3 \wedge x_2) \vee (\bar{x}_1 \wedge x_3) \vee (x_2 \wedge x_3)$ on $[B_2; -, \vee, \wedge]$.

- (a) Simplify this expression using basic Boolean algebra laws.

- (b) Write this expression in minterm normal form.
- (c) Write out the table for the given function defined by f and compare it to the tables of the functions in parts a and b.
- (d) How many possible different functions in three variables on $[B_2; -, \vee, \wedge]$ are there?

B Exercise

3. Let $[B; -, \vee, \wedge]$ be a Boolean algebra of order 4, and let f be a Boolean function of two variables on B .
- (a) How many elements are there in the domain of f ?
- (b) How many different Boolean functions are there of two, variables? Three variables?
- (c) Determine the minterm normal form of $f(x_1, x_2) = x_1 \vee x_2$.
- (d) If $B = \{0, a, b, 1\}$, define a function from B^2 into B that is not a Boolean function.

13.7 A Brief Introduction to the Application of Boolean Algebra to Switching Theory

The algebra of switching theory is Boolean algebra. The standard notation used for Boolean algebra operations in most logic design/switching theory texts is $+$ for \vee and \cdot for \wedge . Complementation is as in this text. Therefore, $(x_1 \wedge \bar{x}_2) \vee (x_1 \wedge x_2) \vee (\bar{x}_1 \wedge x_2)$ becomes $x_1 \cdot \bar{x}_2 + x_1 \cdot x_2 + \bar{x}_1 \cdot x_2$, or simply $x_1 \bar{x}_2 + x_1 x_2 + \bar{x}_1 x_2$. All concepts developed previously for Boolean algebras hold. The only change is purely notational. We make the change in this section solely to introduce the reader to another frequently used notation. Obviously, we could have continued the discussion with our previous notation.

The simplest switching device is the on-off switch. If the switch is closed, on, current will pass through it; if it is open, off, current will not pass through it. If we designate on by true or the logical, or Boolean, 1, and off by false, the logical, or Boolean, 0, we can describe electrical circuits containing switches by logical, or Boolean, expressions. The expression $x_1 \cdot x_2$ represents the situation in which a series of two switches appears in a circuit (see Figure 13.7. 1a). In order for current to flow through the circuit, both switches must be on, that is, have the value 1.

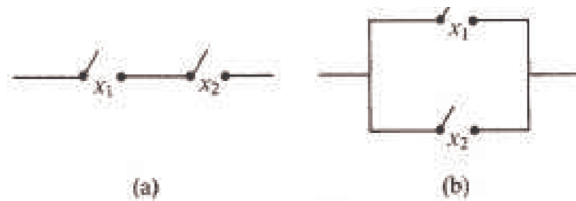
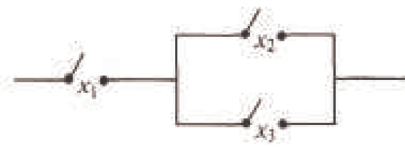


FIGURE 13.7.1

Similarly, a pair of parallel switches, as in Figure 13.7.1b, is described algebraically by $x_1 + x_2$. Many of the concepts in Boolean algebra can be thought of in terms of switching theory. For example, the distributive law in Boolean algebra (in $+$, \cdot notation) is: $x_1 \cdot (x_2 + x_3) = x_1 \cdot x_2 + x_1 \cdot x_3$. Of course, this says the expression on the left is always equivalent to that on the right. The switching circuit analogue of the above statement is that Figure 13.7.2a is equivalent (as an electrical circuit) to Figure 13.7.2b.

The circuits in a digital computer are composed of large quantities of switches that can be represented as in Figure 13.7.2 or can be thought of as boxes or gates with two or more inputs (except for the NOT gate) and one output. These are often drawn as in Figure 13.7.3. For example, the OR gate, as the name implies, is the logical/Boolean OR function. The on-off switch function in Figure 13.7.3a in gate notation is Figure 13.7.3b.

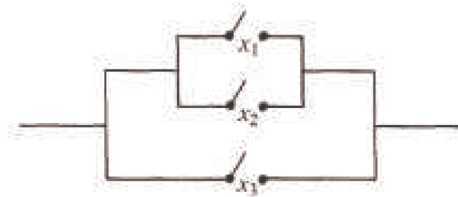


(a)



(b)

FIGURE 13.7.2



(a)



(b)

FIGURE 13.7.3

Either diagram indicates that the circuit will conduct current if and only if $f(x_1, x_2, x_3)$ is true, or 1. We list the gate symbols that are widely used in switching theory in Figure 13.7.4 with their names. The names mean, and are read, exactly as they appear. For example, NAND means "not x_1 and x_2 " or algebraically, $\overline{x_1 \wedge x_2}$, or $\overline{x_1} \cdot \overline{x_2}$.

The circuit in Figure 13.7.5a can be described by gates. To do so, simply keep in mind that the Boolean function $f(x_1, x_2) = x_1 \cdot \overline{x_2}$ of this circuit contains two operations. The operation of complementation takes precedence over that of "and," so we have Figure 13.7.5b.

Example 13.7.1. The switching circuit in Figure 13.7.6a can be expressed through the logic, or gate, circuit in Figure 13.7.6b.







Operation	Symbol		Logical/Boolean Function	
	read	input output	Mathematics notation	Switch Theory notation
AND	and	 $f(x_1, x_2) = x_1 \cdot x_2$	$f(x_1, x_2) = x_1 \wedge x_2$	$f(x_1, x_2) = x_1 \cdot x_2$
OR	or	 $f(x_1, x_2) = x_1 + x_2$	$f(x_1, x_2) = x_1 \vee x_2$	$f(x_1, x_2) = x_1 + x_2$
NOT	not	 $f(x_1) = \overline{x_1}$	$f(x_1) = \overline{x_1}$	$f(x_1) = \overline{x_1}$
NAND	not and	 $f(x_1, x_2) = \overline{x_1 + x_2}$	$f(x_1, x_2) = \overline{x_1 \wedge x_2}$	$f(x_1, x_2) = \overline{x_1 \cdot x_2}$
NOR	not or	 $f(x_1, x_2) = \overline{x_1 + x_2}$	$f(x_1, x_2) = \overline{x_1 \vee x_2}$	$f(x_1, x_2) = \overline{x_1 + x_2}$
Exclusive OR	Exclusive or	 $f(x_1, x_2) = x_1 \oplus x_2$	$f(x_1, x_2) = x_1 \oplus x_2$	$f(x_1, x_2) = x_1 \oplus x_2$

FIGURE 13.7.4

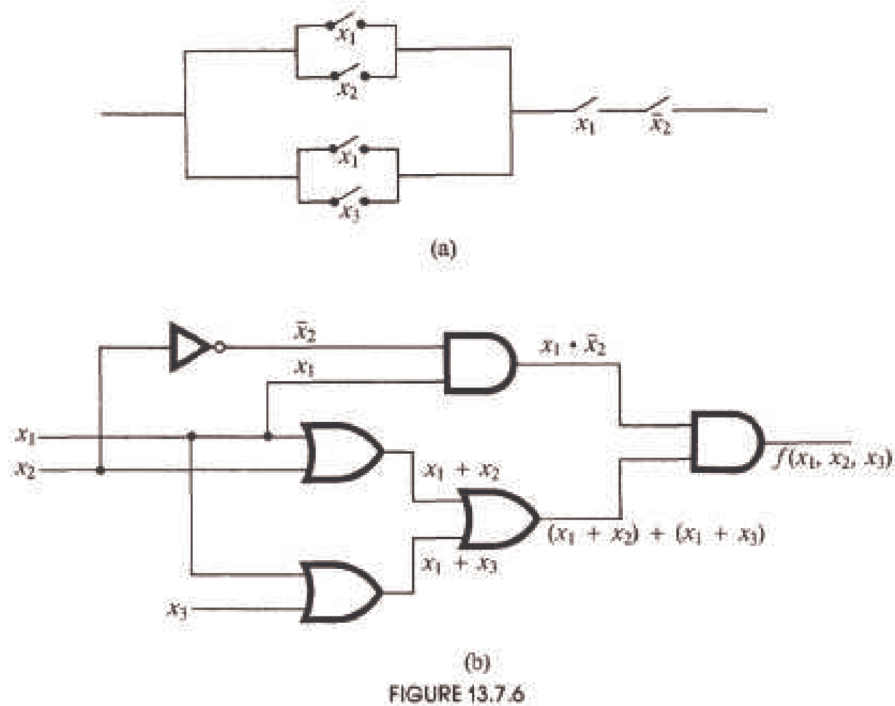
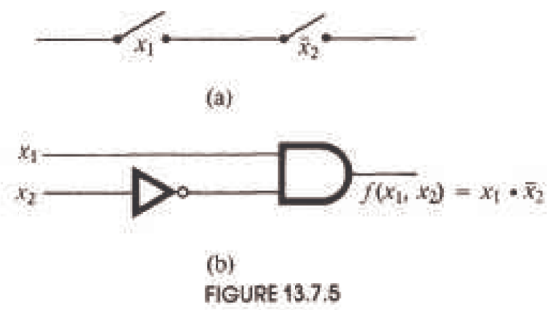
We leave it to the reader to analyze both figures and to convince him- or herself that they do describe the same circuit. The circuit can be described algebraically as

$$f(x_1, x_2, x_3) = ((x_1 + x_2) + (x_1 + x_3)) \cdot x_1 \cdot \overline{x_2}.$$

We can use basic Boolean algebra laws to simplify or minimize this Boolean function (circuit):

$$\begin{aligned}
 f(x_1, x_2, x_3) &= ((x_1 + x_2) + (x_1 + x_3)) \cdot x_1 \cdot \overline{x_2} \\
 &= (x_1 + x_2 + x_3) \cdot x_1 \cdot \overline{x_2} \\
 &= (x_1 \cdot x_1 \cdot \overline{x_2} + x_2 \cdot x_1 \cdot \overline{x_2} + x_3 \cdot x_1 \cdot \overline{x_2}) \\
 &= x_1 \cdot \overline{x_2} + 0 \cdot x_1 + x_3 \cdot x_1 \cdot \overline{x_2} \\
 &= x_1 \cdot \overline{x_2} + x_3 \cdot x_1 \cdot \overline{x_2} \\
 &= x_1 \cdot (\overline{x_2} + \overline{x_2} \cdot x_3) \\
 &= x_1 \cdot \overline{x_2} \cdot (1 + x_3) \\
 &= x_1 \cdot \overline{x_2}.
 \end{aligned}$$

The circuit for f may be described as in Figure 13.7.5. This is a less expensive circuit since it involves considerably less hardware.



The table for f is:

x_1	x_2	x_3	$f(x_1, x_2, x_3)$
0	0	0	0
0	0	1	0
0	1	0	0
0	1	1	0
1	0	0	1
1	0	1	1
1	1	0	0
1	1	1	0

The Venn diagram that represents f is the shaded portion in Figure 13.7.7. From this diagram, we can read off the minterm normal form of f :

$$f(x_1, x_2, x_3) = x_1 \cdot \bar{x}_2 \cdot \bar{x}_3 + x_1 \cdot \bar{x}_2 \cdot x_3.$$

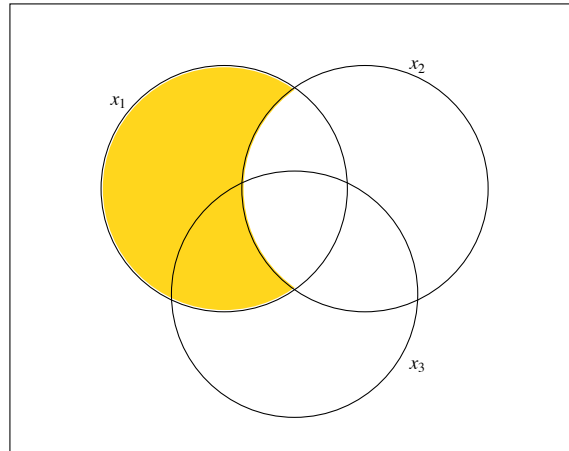


Figure 13.7.7

The circuit (gate) diagram appears in Figure 13.7.8.

How do we interpret this? We see that $f(x_1, x_2, x_3) = 1$ when $x_1 = 1$, $x_2 = 0$, and $x_3 = 0$ or $x_3 = 1$. Current will be conducted through the circuit when switch x_1 is on, switch x_2 is off, and when switch x_3 is either off or on.

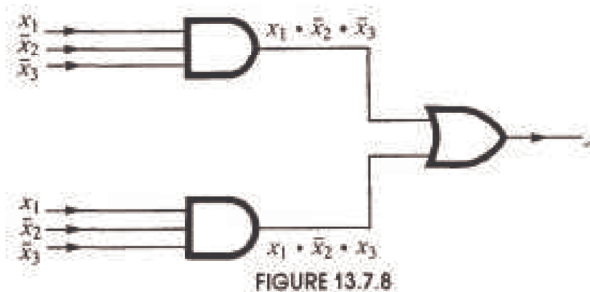


FIGURE 13.7.8

We close this section with a brief discussion of minimization, or reduction, techniques. We have discussed two in this text: algebraic (using basic Boolean rules) reduction and the minterm normal form technique. Other techniques are discussed in switching theory texts. When one reduces a given Boolean function, or circuit, it is possible to obtain a circuit that does not look simpler, but may be more cost effective, and is, therefore, simpler with respect to time. We illustrate with an example.

Example 13.7.2. Consider the Boolean function of Figure 13.7.9a is $f(x_1, x_2, x_3, x_4) = ((x_1 \cdot \overline{x_2}) \cdot \overline{x_3}) \cdot x_4$, which can also be diagrammed as in Figure 13.7.9b.

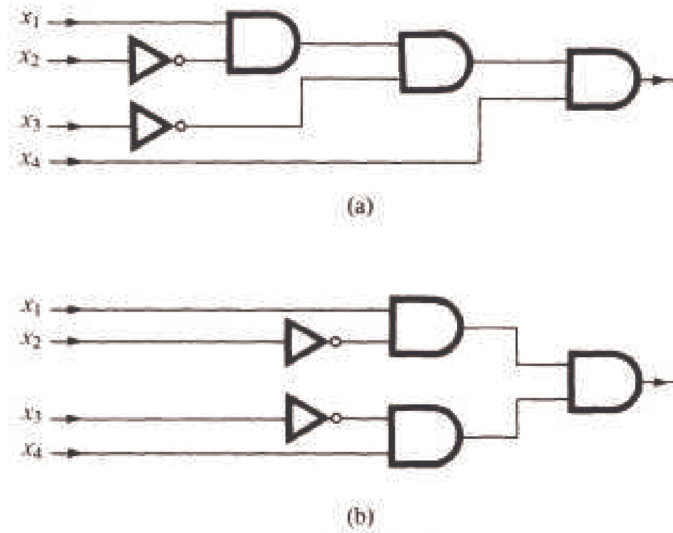


FIGURE 13.7.9

Is Circuit b simpler than Circuit a? Both circuits contain the same number of gates, so the hardware costs (costs per gate) would be the same. Hence, intuitively, we would guess that they are equivalent with respect to simplicity. However, the signals x_3 and x_4 in Circuit a pass through three levels of gating before reaching the output. All signals in Circuit b go through only two levels of gating (disregard the NOT gate when counting levels). Each level of logic (gates) adds to the time delay of the development of a signal at the output. In computers, we want the time delay to be as small as possible. Frequently, speed can be increased by decreasing the number of levels in a circuit. However, this frequently forces a larger number of gates to be used, thus increasing costs. One of the more difficult jobs of a design engineer is to balance off speed with hardware costs (number of gates).

One final remark on notation: The circuit in Figure 13.7.10a can be written as in Figure 13.7.10b, or simply as in Figure 13.7.10c.

EXERCISES FOR SECTION 13.7

A Exercises

1. (a) Write all inputs and outputs from Figure 13.7.11 and show that its Boolean function is $f(x_1, x_2, x_3) = ((x_1 + x_2) \cdot x_3) \cdot (x_1 + x_2)$.
- (b) Simplify f algebraically.
- (c) Find the minterm normal form of f .
- (d) Draw and compare the circuit (gate) diagram of parts b and c above.
- (e) Draw the on-off switching diagram of f in part a.

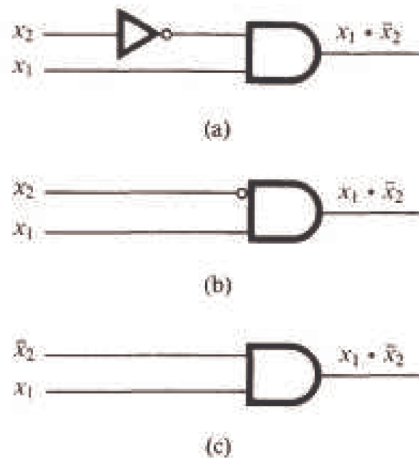


FIGURE 13.7.10

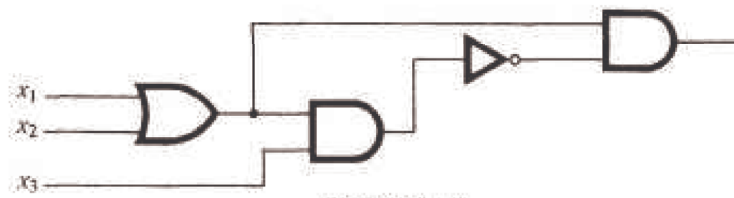


FIGURE 13.7.11

(f) Write the table of the Boolean function f in part a and interpret the results.

2. Given Figure 13.7.12:

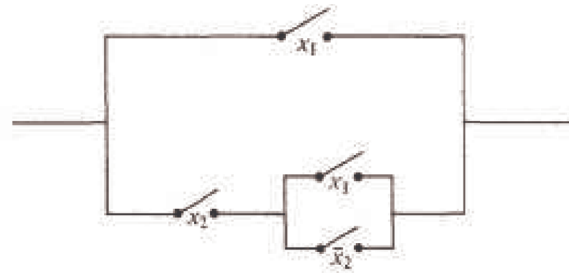


FIGURE 13.7.12

- Write the Boolean function that represents the given on-off circuit.
- Show that the Boolean function obtained in answer to part a can be reduced to $f(x_1, x_2) = x_1$. Draw the on-off circuit diagram of this simplified representation.
- Draw the circuit (gate) diagram of the given on-off circuit diagram.
- Determine the minterm normal of the Boolean function found in the answer to part a or given in part b; they are equivalent.
- Discuss the relative simplicity and advantages of the circuit gate diagrams found in answer to parts c and d.

3. (a) Write the circuit (gate) diagram of

$$f(x_1, x_2, x_3) = (x_1 \bullet x_2 + x_3) \bullet (x_2 + x_3) + x_3.$$

- Simplify the function in part a by using basic Boolean algebra laws.
- Write the circuit (gate) diagram of the result obtained in part b.
- Draw the on-off switch diagrams of parts a and b.

4. Consider the Boolean function

$$f(x_1, x_2, x_3, x_4) = x_1 + (x_2 \bullet (\overline{x_1} + x_4) + x_3 \bullet (\overline{x_2} + \overline{x_4})).$$

- Simplify f algebraically.
- Draw the switching (on-off) circuit of f and the reduction of f obtained in part a.
- Draw the circuit (gate) diagram of f and the reduction of f obtained in answer to part a.

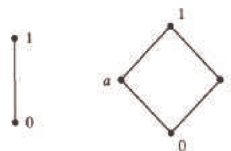
SUPPLEMENTARY EXERCISES FOR CHAPTER 13

Section 13.1

- Draw the Hasse diagram of the relation divides on the set $A = \{1, 2, 3, \dots, 12\}$.
 - For the same set A draw the Hasse diagram for the relation \leq on A .
- For the poset $A = \{1, 2, 3, \dots, 12\}$ under the relation divides find the *lub* and *glb* of the following pairs of numbers if possible: 4 and 6, 2 and 3, 10 and 4, 6 and 9.
 - Repeat part a for the set A , but use the relation \leq .

Section 13.2

- Consider the poset \mathbb{P} under the relation "divides."
 - Compute: $4 \vee 8, 3 \vee 15, 3 \vee 5, 4 \wedge 8, 3 \wedge 15, 3 \wedge 5$ for $[\mathbb{P}, \vee, \wedge]$.
 - Is $[\mathbb{P}, \vee, \wedge]$ a distributive lattice? Explain.
 - Does $[\mathbb{P}, \vee, \wedge]$ have a least element? Does it have a greatest element? If so, what are they?
- Let $[L, \vee, \wedge]$ be a lattice and $a, b \in L$. Prove:
 - $a \vee b = b$ if and only if $a \leq b$.
 - $a \wedge b = a$ if and only if $a \leq b$.
- Let $L = \{0, 1\}$ and define \leq on L by $0 \leq 0 \leq 1 \leq 1$.
 - Draw the Hasse diagram of this poset.
 - Write out the operation table for \vee and \wedge on L observing that they are essentially the standard logical connectives.
 - Define the operations on L^2 componentwise and draw the Hasse diagram for L^2 .
 - Repeat part (c) for L^3 .
- Let $[L_1, \vee, \wedge]$ and $[L_2, \vee, \wedge]$ be lattices. Prove that $[L_1 \times L_2, \vee, \wedge]$ is a lattice when the operations are defined componentwise as we did for algebraic systems in Section 11.6.
 - Let L_1 and L_2 be lattices whose posets have the following Hasse diagrams respectively. List the elements in the lattice $L_1 \times L_2$.



(c) Compute:

$$(0, a) \vee (0, b)$$

$$(0, a) \wedge (0, b)$$

$$(1, a) \vee (1, b)$$

$$(1, a) \wedge (1, b)$$

$$(0, 1) \vee (1, 0)$$

$$\text{and } (0, 1) \wedge (1, 0).$$

Use this information as an aid to draw the Hasse diagram for $L_1 \times L_2$.

- Is $A = \{1, 2, 3, \dots, 12\}$ a lattice under the relation "divides"? Explain.
 - Is the set A above a lattice under the relation "less than or equal to"? Explain.

Section 13.3

- Using the rules of Boolean algebra, reduce the expression $\overline{(x_1 \vee x_2)} \vee (\overline{x_1} \wedge x_2) \vee (x_1 \wedge x_2)$ to the equivalent expression $\overline{x_1} \vee x_2$. Justify each step.

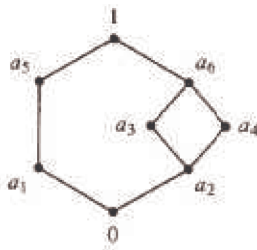
9. Using the rules of Boolean algebra, reduce the expression $(x + y) \cdot (x + \bar{y})$ to a simpler expression.
10. Even a cursory examination of the basic laws for Boolean algebra (Table 13.3.1), for logic (Table 3.4.1), and for sets (Section 4.2) will indicate that they are the same in three different languages: they are isomorphic to one another as Boolean algebras.
- (a) Fill out the following table to illustrate the above concept:

	comparable connectives		
Sets	\cup		
Logic		\wedge	\neg
Boolean Algebra	\leq		

(b) Since the above algebras are isomorphic as Boolean algebras, any theorem true in one is true in the other two. Translate each of the following statements into the language of the other two.

- (i) $p \rightarrow q$ if and only if $\neg q \rightarrow \neg p$.
- (ii) If $A \subseteq B$ and $A \subseteq C$ then $A \subseteq B \cap C$
- (iii) If $a \geq b$ and $a \geq c$ then $a \geq b \vee c$.

11. (a) Determine the complements of each element described by the following Hasse diagram:



(b) Is the above lattice a Boolean algebra?

12. (a) Determine the complement of each element in the lattice D_{50} .
- (b) Is D_{50} a Boolean algebra? Explain.

Section 13.4

13. (a) Use the Theorem 13.4.2 and its Corollaries to determine which of the following are Boolean algebras:
- (a) D_{20} (b) D_{27} (c) D_{35} (d) D_{210}
- (b) Notice that D_n is a Boolean algebra if and only if n is a product of distinct primes. Such an integer is called *square free*. What are the atoms of D_n if n is square free?
14. Let $[B, -, \vee, \wedge]$ be any Boolean algebra of order 8. Find a Boolean algebra of sets that is isomorphic to B . How many atoms must B have?

Section 13.5

15. (a) List all sub-Boolean algebras of order 4 in B_2^3
- (b) How many sub-Boolean algebras of order 4 are there in B_2^n , $n \geq 4$?
- (c) Discuss how the selection of atoms in a sub-Boolean algebra can be used to answer questions such as the one in part (b).
16. Prove that Boolean algebras $B_2^m \times B_2^n$ and B_2^{m+n} are isomorphic.

Section 13.6

17. Find the minterm normal form of the Boolean expression $(\bar{x}_1 \vee x_2) \wedge x_3$
18. Find the minterm normal form of the Boolean expression
- $$x_4 \wedge (x_3 \vee x_2 \vee x_1) \vee x_3 \wedge (x_2 \vee x_1) \vee x_2 \wedge x_1$$
19. Let B be a Boolean algebra of order 2.

(a) How many rows are there in the table of a Boolean function of 3 variables? Of n variables?

(b) How many different Boolean functions of 3 variables and of n variables are there?

20. Let B be a Boolean algebra of order 2.

(a) How many different minterm normal forms are there for Boolean expressions of 2 variables over B ? List them.

(b) How many different minterm normal forms are there for Boolean expressions of 3 variables over B ?

Section 13.7

21. Consider the following Boolean expression:

$$f(x_1, x_2, x_3) = ((x_1 + x_2 + x_3) \cdot \overline{x_1} + x_1 + \overline{x_2}) \cdot x_1 \cdot \overline{x_3}$$

(a) Draw the switching circuit of f .

(b) Draw the gate diagram of f .

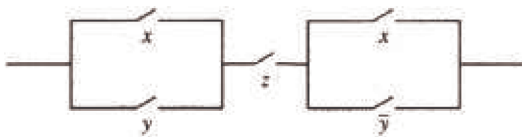
(c) Simplify f algebraically and draw the switching circuit and gate diagrams of this simplified version of f .

22. Assume that each of the three members of a committee votes *yes* or *no* on a proposal by pressing a button that closes a switch for *yes* and does nothing for *no*. Devise as simple a switching-circuit as you can that will allow current to pass when and only when at least two of the members vote in the affirmative.

23. (a) Find the Boolean function of this network:

(b) Draw an equivalent

24. Given the switching circuit

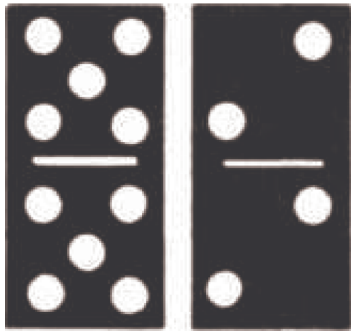


(a) Express the the switching circuit algebraically.

(b) Draw the gate diagram of the expression obtained in part a.

(c) Simplify the expression in part a and draw the switching-circuit and gate diagram for the simplified expression.

chapter 14



Monoids and Automata

GOALS

At first glance, the two topics that we will discuss in this chapter seem totally unrelated. The first is monoid theory, which we touched upon in Chapter 11. The second is automata theory, in which computers and other machines are described in abstract terms. After short independent discussions of these topics, we will describe how the two are related in the sense that each monoid can be viewed as a machine and each machine has a monoid associated with it.

14.1 Monoids

Recall the definition of a monoid:

Definition: Monoid. A monoid is a set M together with a binary operation $*$ with the properties
 (a) $*$ is associative: $(a * b) * c = a * (b * c)$ for all $a, b, c \in M$, and
 (b) $*$ has an identity: there exists $e \in M$ such that for all $a \in M$, $a * e = e * a = a$.

Note: Since the requirements for a group contain the requirements for a monoid, every group is a monoid.

Example 14.1.1.

- (a) The power set of any set together with any one of the operations intersection, union, or symmetric difference is a monoid.
- (b) The set of integers, \mathbb{Z} , with multiplication, is a monoid. With addition, \mathbb{Z} is also a monoid.
- (c) The set of $n \times n$ matrices over the integers, $M_n(\mathbb{Z})$, $n \geq 2$, with matrix multiplication, is a monoid. This follows from the fact that matrix multiplication is associative and has an identity, I_n . This is an example of a noncommutative monoid since there are matrices, A and B , for which $AB \neq BA$.
- (d) $[\mathbb{Z}_n, \times_n]$, $n \geq 2$, is a monoid with identity 1.
- (e) Let X be a nonempty set. The set of all functions from X into X , often denoted X^X , is a monoid over function composition. In Chapter 7, we saw that function composition is associative. The function $i: X \rightarrow X$ defined by $i(a) = a$ is the identity element for this system. This is another example of a noncommutative monoid, provided $|X|$ is greater than 1.

If X is finite, $|X^X| = |X|^{|X|}$. For example, if $B = \{0, 1\}$, $|B^B| = 4$. The functions z , u , i , and t , defined by the graphs in Figure 14.1.1, are the elements of B^B . This monoid is not a group. Do you know why?

One reason that B^B is noncommutative is that $tz \neq zt$, since $(tz)(0) = 1$ and $(zt)(0) = 0$.

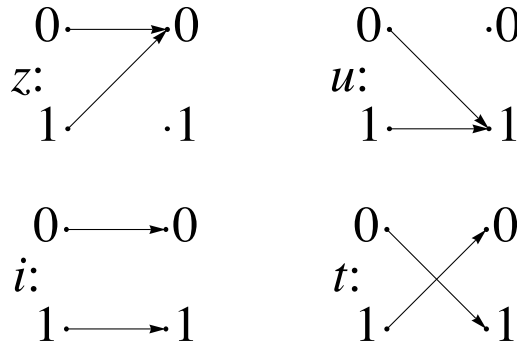


Figure 14.1.1
The four elements of B^B

GENERAL CONCEPTS AND PROPERTIES OF MONOIDS

Virtually all of the group concepts that were discussed in Chapter 11 are applicable to monoids. When we introduced subsystems, we saw that a submonoid of monoid M is a subset of M —that is, it itself is a monoid with the operation of M . To prove that a subset is a submonoid, you can apply the following algorithm.

Theorem/Algorithm 14.1.1. Let $[M; *]$ be a monoid and K is a nonempty subset of M , K is a submonoid of M if and only if:

- (a) If $a, b \in K$, then $a * b \in K$ (i.e., K is closed under $*$), and
- (b) the identity of M belongs to K .

Often we will want to discuss the smallest submonoid that includes a certain subset S of a monoid M . This submonoid can be defined recursively by the following definition.

Definition: Submonoid Generated by a Set. If S is a subset of monoid $[M; *]$, the submonoid generated by S , $\langle S \rangle$, is defined by:

- (a) (Basis) (i) $a \in S \Rightarrow a \in \langle S \rangle$, and (ii) the identity of M belongs to $\langle S \rangle$;
- (b) (Recursion), $a, b \in \langle S \rangle \Rightarrow a * b \in \langle S \rangle$.

Note: If $S = \{a_1, a_2, \dots, a_n\}$, we write $\langle a_1, a_2, \dots, a_n \rangle$ in place of $\langle \{a_1, a_2, \dots, a_n\} \rangle$.

Example 14.1.2.

- (a) In $[\mathbb{Z}; +]$, $\langle 2 \rangle = \{0, 2, 4, 6, 8, \dots\}$.
- (b) The power set of \mathbb{Z} , $\mathcal{P}(\mathbb{Z})$, over union is a monoid with identity \emptyset . If $S = \{\{1\}, \{2\}, \{3\}\}$, then $\langle S \rangle$ is the power set of $\{1, 2, 3\}$. If $S = \{ \{n\} : n \in \mathbb{Z} \}$, then $\langle S \rangle$ is the set of finite subsets of the integers.

MONOID ISOMORPHISMS

Two monoids are *isomorphic* if and only if there exists a translation rule between them so that any true proposition in one monoid is translated to a true proposition in the other.

Example 14.1.3. $M = [\mathcal{P}\{1, 2, 3\}, \cap]$ is isomorphic to $M_2 = [\mathbb{Z}_2^3; \cdot]$, where the operation in M_2 is componentwise *mod 2* multiplication.

A translation rule is that if $A \subseteq \{1, 2, 3\}$, then it is translated to (d_1, d_2, d_3) where $d_i = \begin{cases} 1 & \text{if } i \in A \\ 0 & \text{if } i \notin A \end{cases}$. Two cases of how this translation rule works are:

$$\begin{array}{ccc}
 \{1, 2, 3\} \text{ is the identity for } M_1, & \text{and} & \{1, 2\} \cap \{2, 3\} = \{2\} \\
 \updownarrow & & \updownarrow \quad \updownarrow \quad \updownarrow \quad \updownarrow \\
 (1, 1, 1) \text{ is the identity for } M_2, & \text{and} & (1, 1, 0) \cdot (0, 1, 1) = (0, 1, 0).
 \end{array}$$

A more precise definition of a monoid isomorphism is identical to the definition of a group isomorphism (see Section 11.7).

EXERCISES FOR SECTION 14.1

A Exercises

1. For each of the subsets of the indicated monoid, determine whether the subset is a sub monoid.

- (a) $S_1 = \{0, 2, 4, 6\}$ and $S_2 = \{1, 3, 5, 7\}$ in $[\mathbb{Z}_8; \times_8]$.
- (b) $\{f \in \mathbb{N}^{\mathbb{N}} : f(n) \leq n, \forall n \in \mathbb{N}\}$ and $\{f \in \mathbb{N}^{\mathbb{N}} : f(1) = 2\}$ in $\mathbb{N}^{\mathbb{N}}$.

(c) $\{A \subseteq \mathbb{Z} : A \text{ is finite}\}$ and $\{A \subseteq \mathbb{Z} : A^c \text{ is finite}\}$ in $[\mathcal{P}(\mathbb{Z}); \cup]$.

2. For each subset, describe the submonoid that it generates.

(a) $\{3\}$ and $\{0\}$ in $[\mathbb{Z}_{12}; \times_{12}]$

(b) $\{5\}$ in $[\mathbb{Z}_{25}; \times_{25}]$

(c) the set of prime numbers and $\{2\}$ in $[\mathbb{P}; \cdot]$

(d) $\{3, 5\}$ in $[\mathbb{N}; +]$

B Exercises

3. Definition: Stochastic Matrix. An $n \times n$ matrix of real numbers is called stochastic if and only if each entry is nonnegative and the sum of entries in each column is 1. Prove that the set of stochastic matrices is a monoid over matrix multiplication.

4. Prove Theorem 14.1.1.

14.2 Free Monoids and Languages

In this section, we will introduce the concept of a language. Languages are subsets of a certain type of monoid, the free monoid over an alphabet. After defining a free monoid, we will discuss languages and some of the basic problems relating to them. We will also discuss the common ways in which languages are defined.

Let A be a nonempty set, which we will call an *alphabet*. Our primary interest will be in the case where A is finite; however, A could be infinite for most of the situations that we will describe. The elements of A are called *letters or symbols*. Among the alphabets that we will use are $B = \{0, 1\}$, ASCII = the set of ASCII characters, and PAS = the Pascal character set (whichever one you use).

Definition: Strings over an Alphabet. A string of length n , $n \geq 1$, over A is a sequence of n letters from A : $a_1 a_2 \dots a_n$. The null string, λ , is defined as the string of length zero containing no letters. The set of strings of length n over A is denoted by A^n . The set of all strings over A is denoted A^* .

Notes:

(a) If the length of string s is n , we write $|s| = n$.

(b) The null string is not the same as the empty set, although they are similar in many ways.

(c) $A^* = A^0 \cup A^1 \cup A^2 \cup A^3 \cup \dots$ and if $i \neq j$, $A^i \cap A^j = \emptyset$; that is, $\{A^0, A^1, A^2, A^3, \dots\}$ is a partition of A^* .

(d) An element of A can appear any number of times in a string.

Theorem 14.2.1. If A is countable, then A^* is countable.

Proof: Case 1. Given the alphabet $B = \{0, 1\}$, we can define a bijection from the positive integers into B^* . Each positive integer has a binary expansion $d_k d_{k-1} \dots d_1 d_0$, where each d_j is 0 or 1 and $d_k = 1$. If n has such a binary expansion, then $2^k \leq n < 2^{k+1}$. We define $f: \mathbb{P} \rightarrow B^*$ by $f(n) = f(d_k d_{k-1} \dots d_1 d_0) = d_k \dots d_1 d_0$, where $f(1) = \lambda$. Every one of the 2^k strings of length k are the images of exactly one of the integers between 2^k and $2^{k+1} - 1$. From its definition, f is clearly a bijection; therefore, B^* is countable.

Case 2: A is Finite. We will describe how this case is handled with an example first and then give the general proof. If $A = \{a, b, c, d, e\}$, then we can code the letters in A into strings from B^3 . One of the coding schemes (there are many) is $a \leftrightarrow 000$, $b \leftrightarrow 001$, $c \leftrightarrow 010$, $d \leftrightarrow 011$, and $e \leftrightarrow 100$. Now every string in A^* corresponds to a different string in B^* ; for example, ace would correspond with 000010100 . The cardinality of A^* is equal to the cardinality of the set of strings that can be obtained from this encoding system. The possible coded strings must be countable, since they are a subset of a countable set (B^*); therefore, A^* is countable.

If $|A| = m$, then the letters in A can be coded using a set of fixed-length strings from B^* . If $2^{k-1} < m \leq 2^k$, then there are at least as many strings of length k in B^* as there are letters in A . Now we can associate each letter in A with an element of B^k . Then any string in A^* corresponds to a string in B^* . By the same reasoning as in the example above, A^* is countable.

Case 3: A is Countably Infinite. We will leave this case as an exercise. ■

FREE MONOIDS OVER AN ALPHABET

The set of strings over any alphabet is a monoid under concatenation.

Definition: Concatenation. Let $a = a_1 a_2 \dots a_m$ and $b = b_1 b_2 \dots b_n$ be strings of length m and n , respectively. The concatenation of a with b , $a \langle b$, is the string of length $m + n$: $a_1 a_2 \dots a_m b_1 b_2 \dots b_n$.

Notes:

(a) The null string is the identity element of $[A^*; \text{concatenation}]$. Henceforth, we will denote the monoid of strings over A by A^* .

(b) Concatenation is noncommutative, provided $|A| > 1$.

(c) If $|A_1| = |A_2|$, then the monoids A_1^* and A_2^* are isomorphic. An isomorphism can be defined using any bijection $f: A_1 \rightarrow A_2$. If $a = a_1 a_2 \cdots a_n \in A_1^*$, $f^*(a) = f(a_1) f(a_2) \cdots f(a_n)$ defines a bijection from A_1^* into A_2^* . We will leave it to the reader to convince him or herself that for all $a, b \in A_1^*$, $f^*(a <> b) = f^*(a) <> f^*(b)$.

LANGUAGES

The languages of the world—English, German, Russian, Chinese, and so forth—are called natural languages. In order to communicate in writing in any one of them, you must first know the letters of the alphabet and then know how to combine the letters in meaningful ways. A *formal language* is an abstraction of this situation.

Definition: Formal Language. If A is an alphabet, a formal language over A is a subset of A^* .

Example 14.2.1.

- (a) English can be thought of as a language over the set of letters A, B, \dots, Z (upper and lower case) and other special symbols, such as punctuation marks and the blank. Exactly what subset of the strings over this alphabet defines the English language is difficult to pin down exactly. This is a characteristic of natural languages that we try to avoid with formal languages.
- (b) The set of all ASCII stream files can be defined in terms of a language over ASCII. An ASCII stream file is a sequence of zero or more lines followed by an end-of-file symbol. A line is defined as a sequence of ASCII characters that ends with the two characters CR (carriage return) and LF (line feed). The end-of-file symbol is system-dependent; for example, CTRL/C is a common one.
- (c) The set of all syntactically correct expressions in *Mathematica* is a language over the set of ASCII strings.
- (d) A few languages over B are

$$L_1 = \{s \in B^* \mid s \text{ has exactly as many } 1's \text{ as it has } 0's\},$$

$$L_2 = \{1 <> s <> 0 : s \in B^*\}, \text{ and}$$

$$L_3 = \langle 0, 01 \rangle = \text{the submonoid of } B^* \text{ generated by } \{0, 01\}.$$

TWO FUNDAMENTAL PROBLEMS: RECOGNITION AND GENERATION

The generation and recognition problems are basic to computer programming. Given a language, L , the programmer must know how to write (or generate) a syntactically correct program that solves a problem. On the other hand, the compiler must be written to recognize whether a program contains any syntax errors.

The Recognition Problem: Design an algorithm that determines the truth of $s \in L$ in a finite number of steps for all $a \in A^*$. Any such algorithm is called a *recognition algorithm*.

Definition: Recursive Language. A language is recursive if there exists a recognition algorithm for it.

Example 14.2.2.

- (a) The language of syntactically correct *Mathematica* expressions is recursive.
- (b) The three languages in Example 14.2.1 (d) are all recursive. Recognition algorithms for L_1 and L_2 should be easy for you to imagine. The reason a recognition algorithm for L_3 might not be obvious is that L_3 's definition is more cryptic. It doesn't tell us what belongs to L_3 , just what can be used to create strings in L_3 . This is how many languages are defined. With a second description of L_3 , we can easily design a recognition algorithm. $L_3 = \{s \in B^* : s = \lambda \text{ or } s \text{ starts with a } 0 \text{ and has no consecutive } 1's\}$.

Algorithm 14.2.1: Recognition Algorithm for L_3 . Let $s = s_1 s_2 \cdots s_n \in B^*$. This algorithm determines the truth value of $s \in L_3$. The truth value is returned as the value of Word.

```
(1) Word := true
(2) If  $n > 0$  then
    If  $s_1 = 1$  then Word := false
    else for  $i := 3$  to  $n$ 
        if  $s_{i-1} = 1$  and  $s_i = 1$  then Word := false
```

The Generation Problem. Design an algorithm that generates or produces any string in L . Here we presume that A is either finite or countably infinite; hence, A^* is countable by Theorem 14.2.1, and $L \subseteq A^*$ must be countable. Therefore, the generation of L amounts to creating a list of strings in L . The list may be either finite or infinite, and you must be able to show that every string in L appears somewhere in the list.

Theorem 14.2.2.

- (a) If A is countable, then there exists a generating algorithm for A^* .
- (b) If L is a recursive language over a countable alphabet, then there exists a generating algorithm for L .

Proof:

- (a) Part a follows from the fact that A^* is countable; therefore, there exists a complete list of strings in A^* .
- (b) To generate all strings of L , start with a list of all strings in A^* and an empty list, W , of strings in L . For each string s , use a recognition algorithm (one exists since L is recursive) to determine whether $s \in L$. If s is in L , add it to W ; otherwise "throw it out." Then go to the next string in the list of A^* . ■

Example 14.2.3. Since all of the languages in Example 14.2.2 are recursive, they must have generating algorithms. The one given in the proof of Theorem 14.2.2 is not generally the most efficient. You could probably design more efficient generating algorithms for L_2 and L_3 ; however, a better generating algorithm for L_1 is not quite so obvious.

The recognition and generation problems can vary in difficulty depending on how a language is defined and what sort of algorithms we allow ourselves to use. This is not to say that the means by which a language is defined determines whether it is recursive. It just means that the truth of " L is recursive" may be more difficult to determine with one definition than with another. We will close this section with a discussion of grammars, which are standard forms of definition for a language. When we restrict ourselves to only certain types of algorithms, we can affect our ability to determine whether $s \in L$ is true. In defining a recursive language, we do not restrict ourselves in any way in regard to the type of algorithm that will be used. In Section 14.3, we will consider machines called *finite automata*, which can only perform simple algorithms.

PHRASE STRUCTURE GRAMMARS AND LANGUAGES

One common way of defining a language is by means of a *phrase structure grammar* (or grammar, for short). The set of strings that can be produced using the grammar rules is called the *phrase structure language* (of the grammar).

Example 14.2.4. We can define the set of all strings over B for which all 0s precede all 1s as follows. Define the starting symbol S and establish rules that S can be replaced with any of the following: λ , $0S$, or $S1$. These replacement rules are usually called *production* (or *rewriting*) *rules* and are usually written in the format $S \rightarrow \lambda$, $S \rightarrow 0S$, and $S \rightarrow S1$. Now define L to be the set of all strings that can be produced by starting with S and applying the production rules until S no longer appears. The strings in L are exactly the ones that are described above.

Definition: *Phrase Structure Grammar.* A phrase structure grammar consists of four components:

- (1) A nonempty finite set of terminal characters, T . If the grammar is defining a language over A , T is a subset of A^* .
- (2) A finite set of nonterminal characters, N .
- (3) A starting symbol, $S \in N$.
- (4) A finite set of production rules, each of the form $X \rightarrow Y$, where X and Y are strings over $A \cup N$ such that $X \neq Y$ and X contains at least one nonterminal symbol.

If G is a phrase structure grammar, $L(G)$ is the set of strings that can be obtained by starting with S and applying the production rules a finite number of times until no nonterminal characters remain. If a language can be defined by a phrase structure grammar, then it is called a *phrase structure language*.

Example 14.2.5. The language over B consisting of strings of alternating 0s and 1s is a phrase structure language. It can be defined by the following grammar:

- (1) Terminal characters: λ , 0, and 1,
- (2) Nonterminal characters: S , T , and U ,
- (3) Starting symbol: S ,
- (4) Production rules:

$$\begin{aligned} S &\rightarrow T, S \rightarrow U, S \rightarrow \lambda, S \rightarrow 0, S \rightarrow 1, S \rightarrow 0T, \\ S &\rightarrow 1U, T \rightarrow 10T, T \rightarrow 10, U \rightarrow 01U, U \rightarrow 01 \end{aligned}$$

These rules can be visualized more easily with a graph:

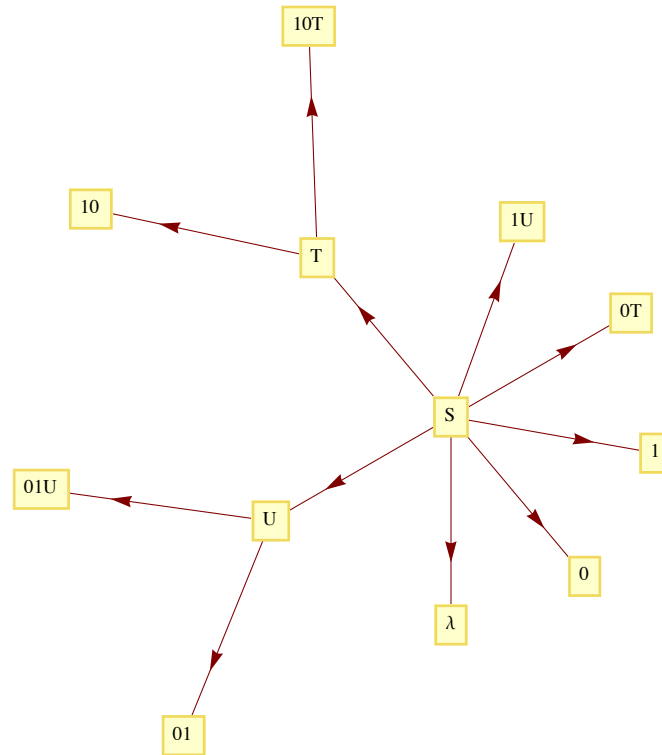


Figure 14.2.1
Production rules for the language of alternating 0's and 1's.

We can verify that a string such as 10101 belongs to the language by starting with S and producing 10101 using the production rules a finite number of times: $S \rightarrow 1U \rightarrow 101U \rightarrow 10101$.

Example 14.2.6. Let G be the grammar with components:

- (1) Terminal symbols = all letters of the alphabet (both upper and lower case) and the digits 0 through 9,
- (2) Nonterminal symbols = $\{I, X\}$,
- (3) Starting symbol: I
- (4) Production rules: $I \rightarrow \alpha$, where α is any letter, $I \rightarrow \alpha X$ for any letter α , $X \rightarrow \beta X$ for any letter or digit β , and $X \rightarrow \beta$ for any letter or digit β .

There are a total of 176 production rules for this grammar. The language $L(G)$ consists of all valid *Mathematica* names.

Backus-Naur form (BNF), A popular alternate form of defining the production rules in a grammar is BNF. If the production rules $A \rightarrow B_1$, $A \rightarrow B_2$, ..., $A \rightarrow B_n$ are part of a grammar, they would be written in BNF as $A ::= B_1 | B_2 | \dots | B_n$. The symbol $|$ in BNF is read as "or," while the $::=$ is read as "is defined as." Additional notations of BNF are that $\{x\}$, represents zero or more repetitions of x and $[y]$ means that y is optional.

Example 14.2.7. A BNF version of the production rules for a *Mathematica* name is

$$\text{letter} ::= a | b | c \dots | z | A | B | \dots | Z$$

$$\text{digit} ::= 0 | 1 | \dots | 9$$

$$I ::= \text{letter} \{ \text{letter} | \text{digit} \}$$

Example 14.2.8. An arithmetic expression can be defined in BNF. For simplicity, we will consider only expressions obtained using addition and multiplication of integers. The terminal symbols are $(,), +, *, -, \text{and the digits } 0 \text{ through } 9$. The nonterminal symbols are E (for expression), T (term), F (factor), and N (number). The starting symbol is E .

$$E ::= E + T | T$$

$$T ::= T * F \mid F$$

$$F ::= (E) \mid N$$

$$N ::= [-] \text{ digit } \{ \text{digit} \}.$$

One particularly simple type of phrase structure grammar is the regular grammar.

Definition: Regular Grammar. A regular (right-hand form) grammar is a grammar whose production rules are all of the form $A \rightarrow t$ and $A \rightarrow tB$, where A and B are nonterminal and t is terminal. A left-hand form grammar allows only $A \rightarrow t$ and $A \rightarrow Bt$. A language that has a regular phrase structure language is called a regular language.

Example 14.2.9.

- (a) The set of *Mathematica* names is a regular language since the grammar by which we defined the set is a regular grammar.
- (b) The language of all strings for which all 0s precede all 1s (Example 14.2.4) is regular; however, the grammar by which we defined this set is not regular. Can you define these strings with a regular grammar?
- (c) The language of arithmetic expressions is not regular.

EXERCISES FOR SECTION 14.2

A Exercises

1. (a) If a computer is being designed to operate with a character set of 350 symbols, how many bits must be reserved for each character? Assume each character will use the same number of bits.
- (b) Do the same for 3,500 symbols.
2. It was pointed out in the text that the null string and the null set are different. The former is a string and the latter is a set, two different kinds of objects. Discuss how the two are similar.
3. What sets of strings are defined by the following grammar?
 - (a) Terminal symbols: λ , 0 and 1
 - (b) Nonterminal symbols: S and E
 - (c) Starting symbol: S
 - (d) Production rules: $S \rightarrow 0S0$, $S \rightarrow 1S1$, $S \rightarrow E$, $E \rightarrow \lambda$, $E \rightarrow 0$, $E \rightarrow 1$.
4. What sets of strings are defined by the following grammar?
 - (a) Terminal symbols: λ , a , b , and c
 - (b) Nonterminal symbols: S , T , U and E
 - (c) Starting symbol: S
 - (d) Production rules: $S \rightarrow aS$, $S \rightarrow T$, $T \rightarrow bT$, $T \rightarrow U$, $U \rightarrow cU$, $U \rightarrow E$, $E \rightarrow \lambda$.
5. Define the following languages over B with phrase structure grammars. Which of these languages are regular?
 - (a) The strings with an odd number of characters.
 - (b) The strings of length 4 or less.
 - (c) The palindromes, strings that are the same backwards as forwards.
6. Define the following languages over B with phrase structure grammars. Which of these languages are regular?
 - (a) The strings with more 0s than 1s.
 - (b) The strings with an even number of 1s.
 - (c) The strings for which all 0s precede all 1s.
7. Prove that if a language over A is recursive, then its complement is also recursive.
8. Use BNF to define the grammars in Exercises 3 and 4.

B Exercise

9. (a) Prove that if X_1, X_2, \dots is a countable sequence of countable sets, the union of these sets, $\bigcup_{i=1}^{\infty} X_i$, is countable.
- (b) Using the fact that the countable union of countable sets is countable, prove that if A is countable, then A^* is countable.

14.3 Automata, Finite-State Machines

In this section, we will introduce the concept of an abstract machine. The machines we will examine will (in theory) be capable of performing many of the tasks associated with digital computers. One such task is solving the recognition problem for a language. We will concentrate on one class of machines, finite-state machines (finite automata). And we will see that they are precisely the machines that are capable of recognizing strings in a regular grammar.

Given an alphabet X , we will imagine a string in X^* to be encoded on a tape that we will call an *input tape*. When we refer to a tape, we might imagine a strip of material that is divided into segments, each of which can contain either a letter or a blank.

The typical abstract machine includes an input device, the *read head*, which is capable of reading the symbol from the segment of the input tape that is currently in the read head. Some more advanced machines have a read/write head that can also write symbols onto the tape. The movement of the input tape after reading a symbol depends on the machine. With a finite-state machine, the next segment of the input tape is always moved into the read head after a symbol has been read. Most machines (including finite-state machines) also have a separate output tape that is written on with a *write head*. The output symbols come from an output alphabet, Z , that may or may not be equal to the input alphabet. The most significant component of an abstract machine is its *memory structure*. This structure can range from a finite number of bits of memory (as in a finite-state machine) to an infinite amount of memory that can be sorted in the form of a tape that can be read from and written on (as in a Turing machine).

Definition: *Finite-State Machine.* A finite-state machine is defined by a quintet (S, X, Z, w, t) where

- (1) $S = \{s_1, s_2, \dots, s_r\}$ is the state set, a finite set that corresponds to the set of memory configurations that the machines can have at any time.
- (2) $X = \{x_1, x_2, \dots, x_m\}$ is the input alphabet.
- (3) $Z = \{z_1, z_2, \dots, z_n\}$ is the output alphabet.
- (4) $w : X \times S \rightarrow Z$ is the output function, which specifies which output symbol $w(x, s) \in Z$ is written onto the output tape when the machine is in state s and the input symbol x is read.
- (5) $t : X \times S \rightarrow S$ is the next-state (or transition) function, which specifies which state $t(x, s) \in S$ the machine should enter when it is in state s and it reads the symbol x .

Example 14.3.1. Many mechanical devices, such as simple vending machines, can be thought of as finite-state machines. For simplicity, assume that a vending machine dispenses packets of gum, spearmint (S), peppermint (P), and bubble (B), for 25¢ each. We can define the input alphabet to be {deposit 25 ¢, press S , press P , press B } and the state set to be {Locked, Select}, where the deposit of a quarter unlocks the release mechanism of the machine and allows you to select a flavor of gum. We will leave it to the reader to imagine what the output alphabet, output function, and next-state function would be. You are also invited to let your imagination run wild and include such features as a coin-return lever and change maker.

Example 14.3.2. The following machine is called a *parity checker*. It recognizes whether or not a string in B^* contains an even number of 1s. The memory structure of this machine reflects the fact that in order to check the parity of a string, we need only keep track of whether an odd or even number of 1s has been detected.

- (1) The input alphabet is $B = \{0, 1\}$.
- (2) The output alphabet is also B .
- (3) The state set is {even, odd}.
- (4, 5) The following table defines the output and next-state functions:

x	s	$w(x, s)$	$t(x, s)$
0	even	0	even
0	odd	1	odd
1	even	1	odd
1	odd	0	even

Note how the value of the most recent output at any time is an indication of the current state of the machine. Therefore, if we start in the even state and read any finite input tape, the last output corresponds to the final state of the parity checker and tells us the parity of the string on the input tape. For example, if the string 11001010 is read from left to right, the output tape, also from left to right, will be 10001100. Since the last character is a 0, we know that the input string has even parity.

An alternate method for defining a finite-state machine is with a transition diagram. A *transition diagram* is a directed graph that contains a node for each state and edges that indicate the transition and output functions. An edge (s_i, s_j) that is labeled x/z indicates that in state s_i the input x results in an output of z and the next state is s_j . That is, $w(x, s_i) = z$ and $t(x, s_i) = s_j$. The transition diagram for the parity checker

appears in Figure 14.3.1. In later examples, we will see that if different inputs, x_i and x_j , while in the same state, result in the same transitions and outputs, we label a single edge $x_i, x_j/z$ instead of drawing two edges with labels x_i/z and x_j/z .

One of the most significant features of a finite-state machine is that it retains no information about its past states that can be accessed by the machine itself. For example, after we input a tape encoded with the symbols 01101010 into the parity checker, the current state will be even, but we have no indication within the machine whether or not it has always been in even state. Note how the output tape is not considered part of the machine's memory. In this case, the output tape does contain a "history" of the parity checker's past states. We assume that the finite-state machine has no way of recovering the output sequence for later use.

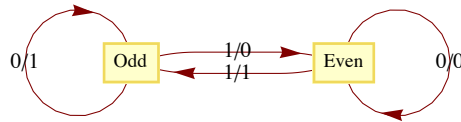


Figure 14.3.1
Transition Diagram for a parity checker

Example 14.3.3. Consider the following simplified version of the game of baseball. To be precise, this machine describes one half-inning of a simplified baseball game. Suppose that in addition to home plate, there is only one base instead of the usual three bases. Also, assume that there are only two outs per inning instead of the usual three. Our input alphabet will consist of the types of hits that the batter could have: out (O), double play (DP), single (S), and home run (HR). The input DP is meant to represent a batted ball that would result in a double play (two outs), if possible. The output alphabet is the numbers 0, 1, and 2 for the number of runs that can be scored as a result of any input. The state set contains the current situation in the inning, the number of outs, and whether a base runner is currently on the base. The list of possible states is then 00 (for 0 outs and 0 runners), 01, 10, 11, and end (when the half-inning is over). The transition diagram for this machine appears in Figure 14.3.2.

Let's concentrate on one state. If the current state is 01, 0 outs and 1 runner on base, each input results in a different combination of output and next-state. If the batter hits the ball poorly (a double play) the output is zero runs and the inning is over (the limit of two outs has been made). A simple out also results in an output of 0 runs and the next state is 11, one out and one runner on base. If the batter hits a single, one run scores (output = 1) while the state remains 01. If a home run is hit, two runs are scored (output = 2) and the next state is 00. If we had allowed three outs per inning, this graph would only be marginally more complicated. The usual game with three bases would be quite a bit more complicated, however.

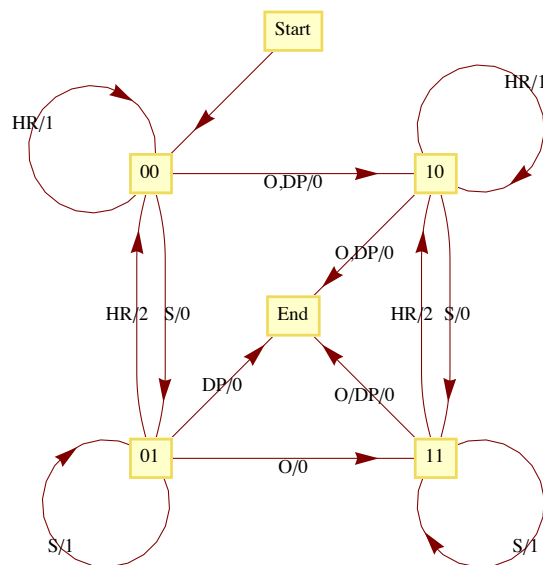


Figure 14.3.2
Transition Diagram for a simplified game of baseball

RECOGNITION IN REGULAR LANGUAGES

As we mentioned at the outset of this section, finite-state machines can recognize strings in a regular language. Consider the language L over $\{a, b, c\}$ that contains the strings of positive length in which each a is followed by b and each b is followed by c . One such string is $bccabc bc$. This language is regular. A grammar for the language would be nonterminal symbols $\{A, B, C\}$ with starting symbol C and production rules $A \rightarrow bB$, cC , $C \rightarrow aA$, $C \rightarrow cC$ and $C \rightarrow c$. A finite-state machine (Figure 14.3.3) that recognizes this language can be constructed with one state for each nonterminal symbol and an additional state (Reject) that is entered if any invalid production takes place. At the end of an input

tape that encodes a string in $\{a, b, c\}^*$, we will know when the string belongs to L based on the final output. If the final output is 1, the string belongs to L and if it is 0, the string does not belong to L . In addition, recognition can be accomplished by examining the final state of the machine. The input string belongs to the language if and only if the final state is C .

The construction of this machine is quite easy: note how each production rule translates into an edge between states other than Reject. For example, $C \rightarrow bB$ indicates that in State C , an input of b places the machine into State B . Not all sets of production rules can be as easily translated to a finite-state machine. Another set of production rules for L is $A \rightarrow aB$, $B \rightarrow bC$, $C \rightarrow cA$, $C \rightarrow cB$, $C \rightarrow cC$ and $C \rightarrow c$. Techniques for constructing finite-state machines from production rules is not our objective here. Hence we will only expect you to experiment with production rules until appropriate ones are found.

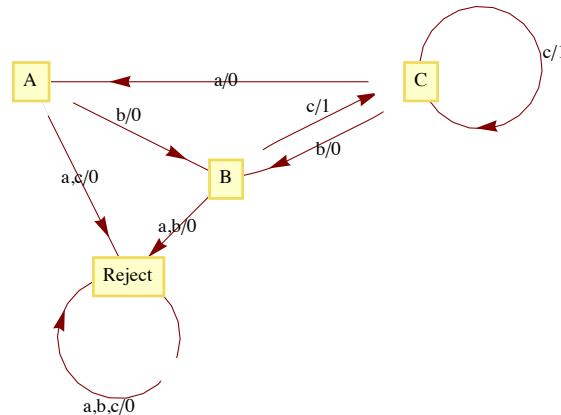
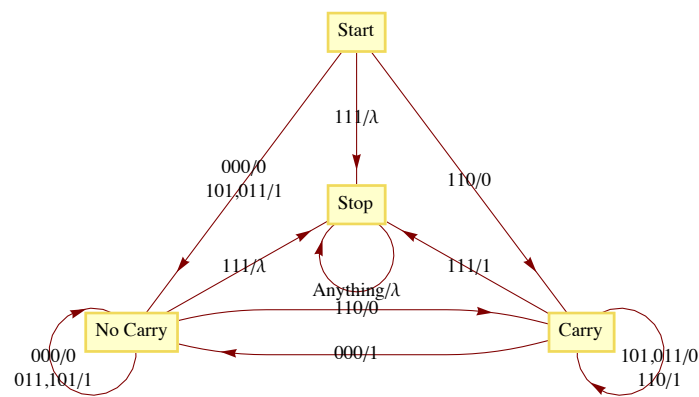


Figure 14.3.3

Example 14.3.4. A finite-state machine can be designed to add positive integers of any size. Given two integers in binary form, $a = a_n a_{n-1} \dots a_1 a_0$ and $b = b_n b_{n-1} \dots b_1 b_0$, the machine will read the input sequence, which is obtained from the digits of a and b reading from right to left,

$$a_0 b_0 (a_0 +_2 b_0), \dots, a_n b_n (a_n +_2 b_n),$$

followed by the special input 111. Note how all possible inputs except the last one must even parity (contain an even number of ones). The output sequence is the sum of a and b , starting with the units digit, and comes from the set $\{0, 1, \lambda\}$. The transition diagram for this machine appears in Figure 14.3.4.

Figure 14.3.4
Transition Diagram for a binary adder

EXERCISES FOR SECTION 14.3

A Exercises

1. Draw a transition diagram for the vending machine described in Example 14.3.1.
2. Construct finite-state machines that recognize the regular languages that you identified in Section 14.2.

3. What is the input set for the machine in Example 14.3.4?
4. What input sequence would be used to compute the sum of 1101 and 0111 (binary integers)? What would the output sequence be?

B Exercise

5. *The Gray Code Decoder.* The finite-state machine defined by the following figure has an interesting connection with the Gray Code (Section 9.4).

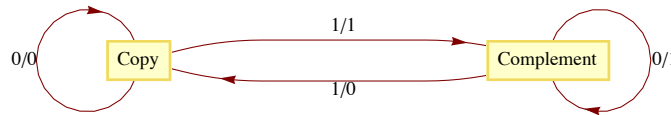


Figure 14.3.5
Gray Code Decoder

Given a string $x = x_1 x_2 \cdots x_n \in B^n$, we may ask where x appears in G_n . Starting in Copy state, the input string x will result in an output string $z \in B^n$, which is the binary form of the position of x in G_n . Positions are numbered from 0 to $2^n - 1$.

- (a) In what positions (0 – 31) do 10110, 00100, and 11111 appear in G_5 ?
- (b) Prove that the Gray Code Decoder always works.

14.4 The Monoid of a Finite-State Machine

In this section, we will see how every finite-state machine has a monoid associated with it. For any finite-state machine, the elements of its associated monoid correspond to certain input sequences. Because only a finite number of combinations of states and inputs is possible for a finite-state machine there is only a finite number of input sequences that summarize the machine. This idea is illustrated best with a few examples.

Example 14.4.1. Consider the parity checker. The following table summarizes the effect on the parity checker of strings in B^1 and B^2 . The row labeled "Even" contains the final state and final output as a result of each input string in B^1 and B^2 when the machine starts in the even state. Similarly, the row labeled "Odd" contains the same information for input sequences when the machine starts in the odd state.

Input String	0	1	00	01	10	11
Even	(Even, 0)	(Odd, 1)	(Even, 0)	(Odd, 1)	(Odd, 1)	(Even, 0)
Odd	(Odd, 1)	(Even, 1)	(Odd, 1)	(Even, 1)	(Even, 0)	(Odd, 1)
Same Effect as			0	1	1	0

Note how, as indicated in the last row, the strings in B^2 have the same effect as certain strings in B^1 . For this reason, we can summarize the machine in terms of how it is affected by strings of length 1. The actual monoid that we will now describe consists of a set of functions, and the operation on the functions will be based on the concatenation operation.

Let T_0 be the final effect (state and output) on the parity checker of the input 0. Similarly, T_1 is defined as the final effect on the parity checker of the input 1. More precisely,

$$T_0(\text{even}) = (\text{even}, 0) \quad \text{and} \quad T_0(\text{odd}) = (\text{odd}, 1),$$

while

$$T_1(\text{even}) = (\text{odd}, 1) \quad \text{and} \quad T_1(\text{odd}) = (\text{even}, 0).$$

In general, we define the operation on a set of such functions as follows: if s, t are input sequences and T_s and T_t are functions as above, then $T_s * T_t = T_{st}$, that is, the result of the function that summarizes the effect on the machine by the concatenation of s with t . Since, for example, 01 has the same effect on the parity checker as 1, $T_0 * T_1 = T_{01} = T_1$. We don't stop our calculation at T_{01} because we want to use the shortest string of inputs to describe the final result. A complete table for the monoid of the parity checker is

*	T_0	T_1
T_0	T_0	T_1
T_1	T_1	T_0

What is the identity of this monoid? The monoid of the parity checker is isomorphic to the monoid $[\mathbb{Z}_2, +_2]$.

This operation may remind you of the composition operation on functions, but there are two principal differences. The domain of T_s is not the codomain of T_t and the functions are read from left to right unlike in composition, where they are normally read from right to left.

You may have noticed that the output of the parity checker echoes the state of the machine and that we could have looked only at the effect on the machine as the final state. The following example has the same property, hence we will only consider the final state.

Example 14.4.2. The transition diagram for the machine that recognizes strings in B^* that have no consecutive 1's appears in Figure 14.4.1. Note how it is similar to the graph in Figure 9.1.1. Only a "reject state" has been added, for the case when an input of 1 occurs while in State a . We construct a similar table to the one in the previous example to study the effect of certain strings on this machine. This time, we must include strings of length 3 before we recognize that no "new effects" can be found.

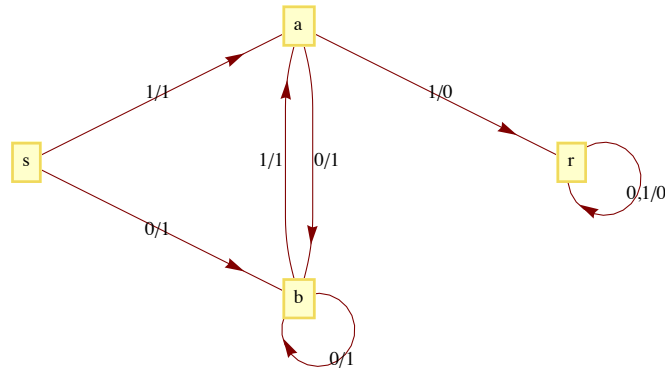


Figure 14.4.1

Inputs	0	1	00	01	10	11	000	001	010	011	100	101	110	111
s	b	a	b	a	b	r	b	a	b	r	b	a	r	r
a	b	r	b	a	r	r	b	a	b	r	r	r	r	r
b	b	a	b	a	b	r	b	a	b	r	b	a	r	r
r	r	r	r	r	r	r	r	r	r	r	r	r	r	r
Same as	0						0	01	0	11	10	1	11	11

The following table summarizes how combinations of the strings 0, 1, 01, 10, and 11 affect this machine.

*	T_0	T_1	T_{01}	T_{10}	T_{11}
T_0	T_0	T_1	T_{01}	T_{10}	T_{11}
T_1	T_{10}	T_{11}	T_1	T_{11}	T_{11}
T_{01}	T_0	T_{11}	T_{01}	T_{11}	T_{11}
T_{10}	T_{10}	T_1	T_1	T_{10}	T_{11}
T_{11}	T_{11}	T_{11}	T_{11}	T_{11}	T_{11}

All the results in this table can be obtained using the previous table. For example,

$$T_{10} * T_{01} = T_{1001} = T_{100} * T_1 = T_{10} * T_1 = T_{101} = T_1$$

and

$$T_{01} * T_{01} = T_{0101} = T_{010} T_1 = T_0 T_1 = T_{01}.$$

Note that none of the elements that we have listed in this table serves as the identity for our operation. This problem can always be remedied by including the function that corresponds to the input of the null string, T_λ . Since the null string is the identity for concatenation of strings, $T_s T_\lambda = T_\lambda T_s = T_s$ for all input strings s .

Example 14.4.3. A finite-state machine called the unit-time delay machine does not echo its current state, but prints its previous state. For this reason, when we find the monoid of the unit-time delay machine, we must consider both state and output. The transition diagram of this machine appears in Figure 14.4.2.

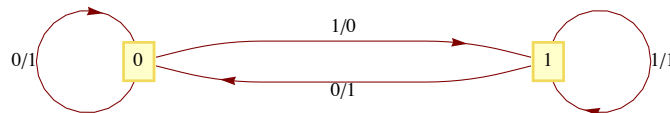


Figure 14.4.2

Input	0	1	00	01	10	11	100 or 000	101 or 001	110 or 101	111 or 011
0	(0, 0)	(1, 0)	(0, 0)	(1, 0)	(0, 1)	(1, 1)	(0, 0)	(1, 0)	(0, 1)	(1, 1)
1	(0, 1)	(1, 1)	(0, 0)	(1, 0)	(0, 1)	(1, 1)	(0, 0)	(1, 0)	(0, 1)	(1, 1)
Same as							00	01	10	11

Again, since no new outcomes were obtained from strings of length 3, only strings of length 2 or less contribute to the monoid of the machine. The table for the strings of positive length shows that we must add T_λ to obtain a monoid.

*	T_0	T_1	T_{00}	T_{01}	T_{10}	T_{11}
T_0	T_{00}	T_{01}	T_{00}	T_{01}	T_{10}	T_{11}
T_1	T_{10}	T_{11}	T_{00}	T_{01}	T_{10}	T_{11}
T_{00}	T_{00}	T_{01}	T_{00}	T_{01}	T_{10}	T_{11}
T_{01}	T_{10}	T_{11}	T_{00}	T_{01}	T_{10}	T_{11}
T_{10}	T_{00}	T_{01}	T_{00}	T_{01}	T_{10}	T_{11}
T_{11}	T_{10}	T_{11}	T_{00}	T_{01}	T_{10}	T_{11}

EXERCISES FOR SECTION 14.4

A Exercise

1. For each of the transition diagrams in Figure 14.4.3, write out tables for their associated monoids. Identify the identity in terms of a string of positive length, if possible. (*Hint*: Where the output echoes the current state, the output can be ignored.)

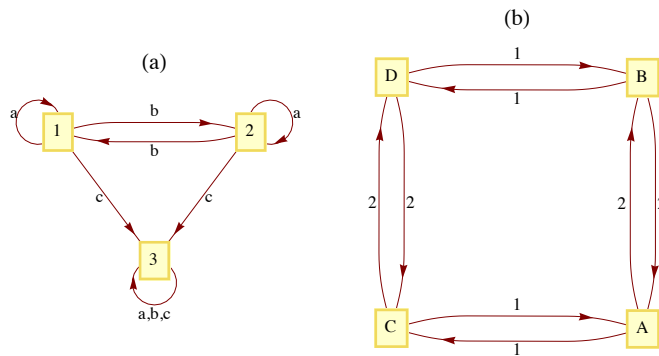


Figure 14.4.3

B Exercise

2. What common monoids are isomorphic to the monoids obtained in the previous exercise?

C Exercise

3. Can two finite-state machines with nonisomorphic transition diagrams have isomorphic monoids?

14.5 The Machine of a Monoid

Any finite monoid $[M, *]$ can be represented in the form of a finite-state machine with input and state sets equal to M . The output of the machine will be ignored here, since it would echo the current state of the machine. Machines of this type are called *state machines*. It can be shown that whatever can be done with a finite-state machine can be done with a state machine; however, there is a trade-off. Usually, state machines that perform a specific function are more complex than general finite-state machines.

Definition: Machine of a Monoid. If $[M, *]$ is a finite monoid, then the machine of M , denoted $m(M)$, is the state machine with state set M , input set M , and next-state function $t : M \times M \rightarrow M$ defined by $t(s, x) = s * x$.

Example 14.5.1. We will construct the machine of the monoid $[\mathbb{Z}_2; +_2]$. As mentioned above, the state set and the input set are both \mathbb{Z}_2 . The next state function is defined by $t(s, x) = s +_2 x$. The transition diagram for $m(\mathbb{Z}_2)$ appears in Figure 14.5.1. Note how it is identical to the transition diagram of the parity checker, which has an associated monoid that was isomorphic to $[\mathbb{Z}_2, +_2]$.

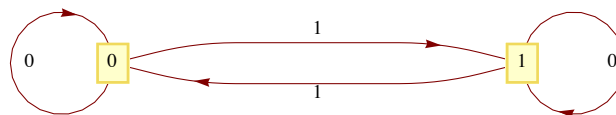


Figure 14.5.1

Example 14.5.2. The transition diagram of the monoids $[\mathbb{Z}_2, \times_2]$ and $[\mathbb{Z}_3, \times_3]$ appear in Figure 14.5.2.

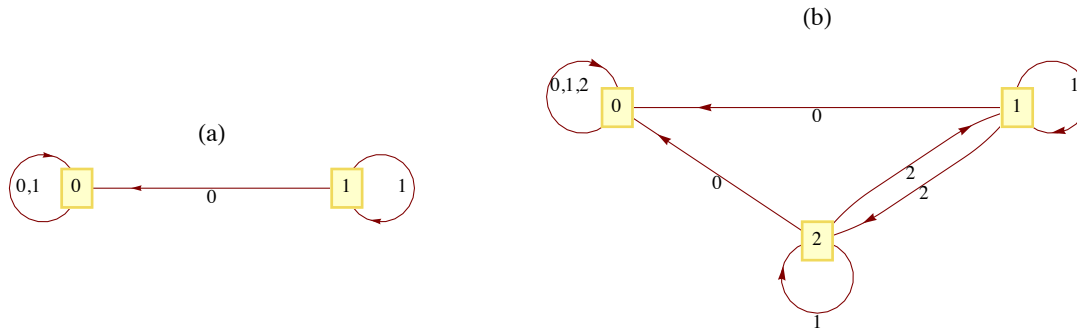


Figure 14.5.2

Example 14.5.3. Let U be the monoid that we obtained from the unit-time delay machine (Example 14.4.3). We have seen that the machine of the monoid of the parity checker is essentially the parity checker. Will we obtain a unit-time delay machine when we construct the machine of U ? We can't expect to get exactly the same machine because the unit-time delay machine is not a state machine and the machine of a monoid is a state machine. However, we will see that our new machine is capable of telling us what input was received in the previous time period. The operation table for the monoid serves as a table to define the transition function for the machine. The row headings are the state values, while the column headings are the inputs. If we were to draw a transition diagram with all possible inputs, the diagram would be too difficult to read. Since U is generated by the two elements, T_0 and T_1 , we will include only those inputs. Suppose that we wanted to read the transition function for the input T_{01} . Since $T_{01} = T_0 T_1$, in any state s , $t(s, T_{01}) = t(t(s, T_0), T_1)$. The transition diagram appears in Figure 14.5.3.

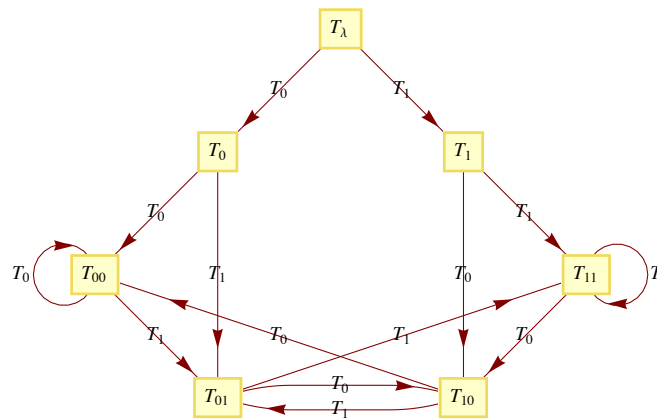


Figure 14.5.3

If we start reading a string of 0s and 1s while in state T_λ and are in state T_{ab} at any one time, the input from the previous time period (not the input that sent us into T_{ab} , the one before that) is a . In states T_λ , T_0 and T_1 , no previous input exists.

EXERCISES FOR SECTION 14.5

A Exercise

1. Draw the transition diagrams for the machines of the following monoids:

(a) $[\mathbb{Z}_4; +_4]$

(b) The direct product of $[\mathbb{Z}_2; \times_2]$ with itself.

B Exercise

2. Even though a monoid may be infinite, we can visualize it as an infinite-state machine provided that it is generated by a finite number of elements. For example, the monoid B^* is generated by 0 and 1. A section of its transition diagram can be obtained by allowing input only from the generating set (Figure 14.5.4a). The monoid of integers under addition is generated by the set $\{-1, 1\}$. The transition diagram for this monoid can be visualized by drawing a small portion of it, as in Figure 14.5.4b.

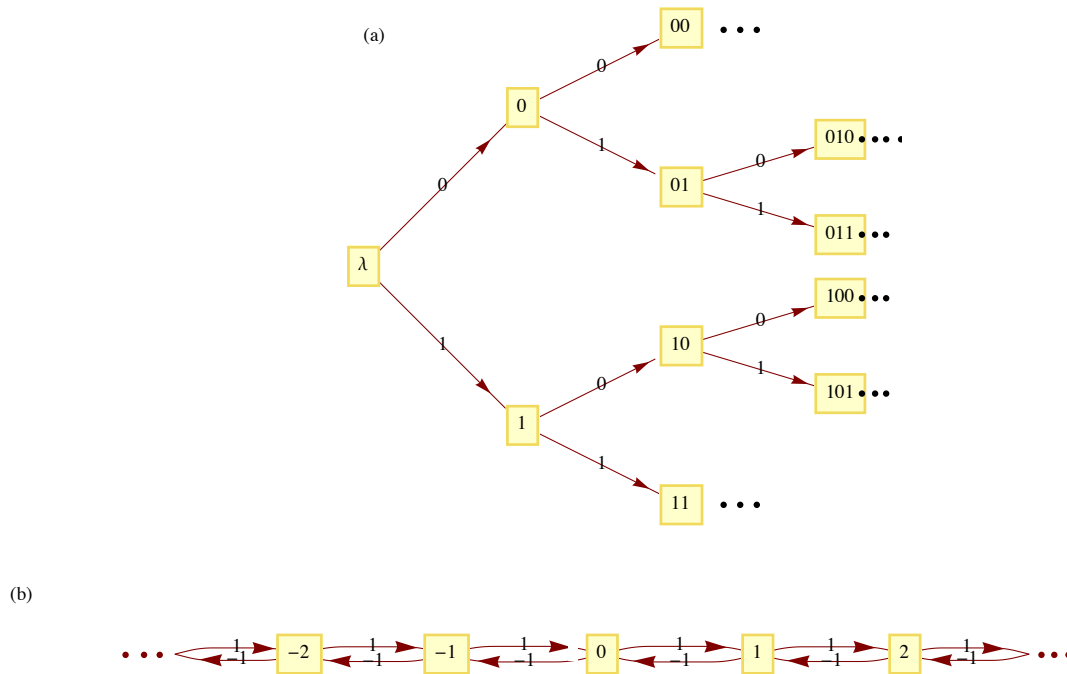


Figure 14.5.4

- (a) Draw a transition diagram for $\{a, b, c\}^*$.
- (b) Draw a transition diagram for $[\mathbb{Z} \times \mathbb{Z}, \text{componentwise addition}]$.

SUPPLEMENTARY EXERCISES FOR CHAPTER 14

Section 14.1

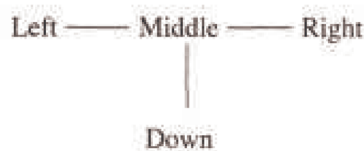
1. Let B be a Boolean algebra and M the set of all Boolean functions on B . Let $*$ be defined on M by $(f * g)(a) = f(a) \wedge g(a)$. Prove that $[M, *]$ is a monoid. Construct the operation table of $[M, *]$ for the case of $B = B_2$.
2. A semigroup is an algebraic system $[S, *]$ with the only axiom that $*$ be associative on S . Prove that if S is a finite set, then there must exist an idempotent element, that is, an $a \in S$ such that $a * a = a$.

Section 14.2

3. What language does the following grammar define? The start symbol is S , the terminal symbols are a and b , and the nonterminal symbols are S and B . The production rules are $S \rightarrow a, S \rightarrow bB, B \rightarrow b, B \rightarrow bS$.
4. What language does the following grammar define? Start symbol = S . nonterminal symbols: T, R . Production rules: $S \rightarrow T, S \rightarrow R, T \rightarrow bR, R \rightarrow aT, T \rightarrow b, R \rightarrow a$.
5. Write a regular grammar for the language L over the alphabet $\{a, b\}$ where L is the set of all strings with exactly two b 's.
6. Write a regular grammar to describe the strings of 0's and 1's that consist of a positive number of 0's surrounded by single 1's. For example, 100001 is one such string.

Section 14.3

7. Draw a finite-state machine to recognize the language in Exercise 5. Have the last output be 1 if the input word is in L , and 0 if it is not in L .
8. Draw a transition diagram for a finite-state machine that recognizes strings in the language of Exercise 6.
9. A finite-state machine moves once every time unit between one of four states called Right, Middle, Left, and Down. The input alphabet is $X = \{00, 01, 10, 11\}$ and the output alphabet is $Y = \{1, 0\}$.



(i) If the machine is in the Middle, Right, or Left, 00 means that it stays where it is; 01 means that it moves one state to the right (e.g. Left to Middle.)—if it can't move any farther to the right, it stays where it is;

10 means that it moves one state to the left.

(ii) Input of 11 means that the machine stays where it is except if it is in the Middle, in which case it enters the Down state.

(iii) If the machine is in the Down state, it stays in that state forever.

(iv) Output is 1 if the state of the machine changes, 0 otherwise.

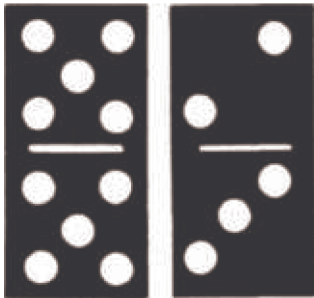
(a) Construct the transition diagram for this finite-state machine.

(b) If $s(0) = \text{Middle}$ and $s(3) = \text{Down}$, list the possible output sequences that could have occurred for $t = 0, 1, 2$.

Section 14.4

10. Write out the operation table for the monoid of the machine in Exercise 9. Section 14.5
11. Draw the transition diagram of the machine of $[M, *]$ in Exercise 1 of these supplementary exercises.

Chapter 15



GROUP THEORY AND APPLICATIONS

GOALS

In Chapter 11, Algebraic Systems, groups were introduced as a typical algebraic system. The associated concepts of subgroup, group isomorphism, and direct products of groups were also introduced. Groups were chosen for that chapter because they are among the simplest types of algebraic systems. Despite this simplicity, group theory abounds with interesting applications, many of which are of interest to the computer scientist. In this chapter we intend to present the remaining important concepts in elementary group theory and some of their applications.

15.1 Cyclic Groups

Groups are classified according to their size and structure. A group's structure is revealed by a study of its subgroups and other properties (e.g., whether it is abelian) that might give an overview of it. Cyclic groups have the simplest structure of all groups.

Definitions: Cyclic Group, Generator. Group G is cyclic if there exists $a \in G$ such that the cyclic subgroup generated by a , $\langle a \rangle$, equals all of G . That is, $G = \{na \mid n \in \mathbb{Z}\}$, in which case a is called a generator of G . The reader should note that additive notation is used for G .

Example 15.1.1. $\mathbb{Z}_{12} = [\mathbb{Z}_{12}, +_{12}]$, where $+_{12}$ is addition modulo 12, is a cyclic group. To verify this statement, all we need to do is demonstrate that some element of \mathbb{Z}_{12} is a generator. One such element is 5; that is, $\langle 5 \rangle = \mathbb{Z}_{12}$. One more obvious generator is 1. In fact, 1 is a generator of every $[\mathbb{Z}_n; +_n]$. The reader is asked to prove that if an element is a generator, then its inverse is also a generator. Thus, $-5 = 7$ and $-1 = 11$ are the other generators of \mathbb{Z}_{12} .

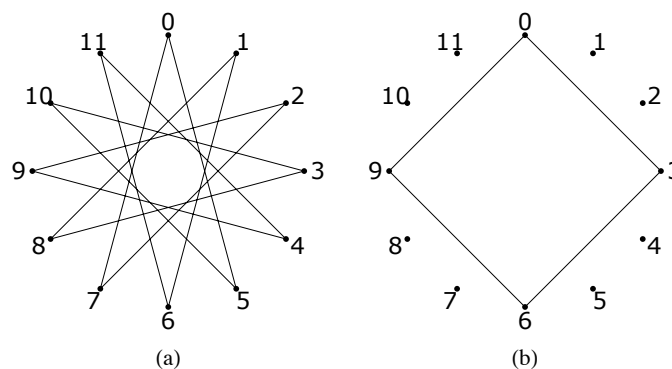


Figure 15.1.1
Examples of "string art"

Figure 15.1.1(a) is an example of "string art" that illustrates how 5 generates \mathbb{Z}_{12} . Twelve tacks are placed along a circle and numbered. A string is tied to tack 0, and is then looped around every fifth tack. As a result, the numbers of the tacks that are reached are exactly the ordered

multiples of 5 modulo 12: 5, 10, 3, ..., 7, 0. Note that if every seventh tack were used, the same artwork would be obtained. If every third tack were connected, as in Figure 15.1.1(b), the resulting loop would only use four tacks; thus 3 does not generate \mathbb{Z}_{12} .

Example 15.1.2. The group of additive integers, $[\mathbb{Z}; +]$, is cyclic:

$$\mathbb{Z} = \langle 1 \rangle = \{n \cdot 1 \mid n \in \mathbb{Z}\}.$$

This observation does not mean that every integer is the product of an integer times 1. It means that

$$\mathbb{Z} = \{0\} \cup \left\{ \overbrace{1 + 1 + \cdots + 1}^{n \text{ terms}} \mid n \in \mathbb{P} \right\} \cup \left\{ \overbrace{(-1) + (-1) + \cdots + (-1)}^{n \text{ terms}} \mid n \in \mathbb{P} \right\}$$

Theorem 15.1.1. If $[G, *]$ is cyclic, then it is abelian.

Proof: Let a be any generator of G and let $b, c \in G$. By the definition of the generator of a group, there exists integers m and n such that $b = ma$ and $c = na$. Thus

$$\begin{aligned} b * c &= (ma) * (na) \\ &= (m + n)a \quad \text{by Theorem 11.3.7(ii)} \\ &= (n + m)a \\ &= (na) * (mb) \\ &= c * b \quad \blacksquare \end{aligned}$$

One of the first steps in proving a property of cyclic groups is to use the fact that there exists a generator. Then every element of the group can be expressed as some multiple of the generator. Take special note of how this is used in theorems of this section.

Up to now we have used only additive notation to discuss cyclic groups. Theorem 15.1.1 actually justifies this practice since it is customary to use additive notation when discussing abelian groups. Of course, some concrete groups for which we employ multiplicative notation are cyclic. If one of its elements, a , is a generator,

$$\langle a \rangle = \{a^n \mid n \in \mathbb{Z}\}$$

Example 15.1.3. The group of positive integers modulo 11 with modulo 11 multiplication, $[\mathbb{Z}_{11}^*; \times_{11}]$, is cyclic. One of its generators is 6: $6^1 = 6, 6^2 = 3, 6^3 = 7, \dots, 6^9 = 2$, and $6^{10} = 1$, the identity of the group.

Example 15.1.4. The real numbers with addition, $[\mathbb{R}; +]$ is a noncyclic group. The proof of this statement requires a bit more generality since we are saying that for all $r \in \mathbb{R}$, (r) is a proper subset of \mathbb{R} . If r is nonzero, the multiples of r are distributed over the real line, as in Figure 15.1.2. It is clear then that there are many real numbers, like $r/2$, that are not in (r) .

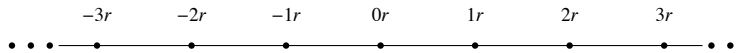


Figure 15.1.2
Elements of (r) , $r > 0$

The following theorem shows that a cyclic group can never be very complicated.

Theorem 15.1.2. If G is a cyclic group, then G is either finite or countably infinite. If G is finite and $|G| = n$, it is isomorphic to $[\mathbb{Z}_n, +_n]$. If G is infinite, it is isomorphic to $[\mathbb{Z}, +]$.

Proof: Case 1: $|G| < \infty$. If a is a generator of G and $|G| = n$, define $\phi: \mathbb{Z}_n \rightarrow G$ by

$$\phi(k) = ka \quad \text{for all } k \in \mathbb{Z}_n$$

Since (a) is finite, we can use the fact that the elements of (a) are the first n nonnegative multiples of a . From this observation, we see that ϕ is a surjection. A surjection between finite sets of the same cardinality must be a bijection. Finally, if $p, q \in \mathbb{Z}_n$,

$$\begin{aligned} \phi(p) + \phi(q) &= pa + qa \\ &= (p + q)a \\ &= (p +_n q)a \quad \text{see exercise 10} \\ &= \phi(p +_n q) \end{aligned}$$

Therefore ϕ is an isomorphism.

Case 2; $|G| = \infty$. We will leave this case as an exercise. \blacksquare

The proof of the next theorem makes use of the division property for integers, which was introduced in Section 11.4: If m, n are integers, $m > 0$, there exist unique integers q (quotient) and r (remainder) such that $n = qm + r$ and $0 \leq r < m$.

Theorem 15.1.3. Every subgroup of a cyclic group is cyclic.

Proof: Let G be cyclic with generator a and let $H \leq G$. If $H = \{e\}$, H has e as a generator. We may now assume that $|H| \geq 2$ and $a \neq e$. Let m be the least positive integer such that ma belongs to H . (This is the key step. It lets us get our hands on a generator of H .) We will now show that $c = ma$ generates H . Suppose that $(c) \neq H$. Then there exists $b \in H$ such that $b \notin (c)$. Now, since b is in G , there exists $n \in \mathbb{Z}$ such that $b = na$. We now apply the division property and divide n by m .

$$b = na = (qm + r)a = (qm)a + ra,$$

where $0 \leq r < m$. We note that r cannot be zero for otherwise we would have $b = na = q(ma) = qc \in (c)$. Therefore,

$$ra = na - (qm)a \in H$$

This contradicts our choice of m because $0 < r < m$. ■

Example 15.1.5. The only proper subgroups of \mathbb{Z}_{10} are $H_1 = \{0, 5\}$ and $H_2 = \{0, 2, 4, 6, 8\}$. They are both cyclic: $H_1 = \langle 5 \rangle$, while $H_2 = \langle 2 \rangle = \langle 4 \rangle = \langle 6 \rangle = \langle 8 \rangle$. The generators of \mathbb{Z}_{10} are 1, 3, 7, and 9.

Example 15.1.6. With the exception of $\{0\}$, all subgroups of \mathbb{Z} are isomorphic to \mathbb{Z} . If $H \leq \mathbb{Z}$, then H is the cyclic subgroup generated by the least positive element of H . It is infinite and so by theorem 15.1.2 it is isomorphic to \mathbb{Z} .

We now cite a useful theorem for computing the order of cyclic subgroups of a cyclic group:

Theorem 15.1.4. *If G is a cyclic group of order n and a is a generator of G , the order of ka is n/d , where d is the greatest common divisor of n and k .*

The proof of this theorem is left to the reader.

Example 15.1.7. To compute the order of (18) in \mathbb{Z}_{30} , we first observe that 1 is a generator of \mathbb{Z}_{30} and $18 = 18(1)$. The greatest common divisor of 18 and 30 is 6. Hence, the order of (18) is $30/6$, or 5.

APPLICATION: FAST ADDERS

At this point, we will introduce the idea of a fast adder, a relatively modern application (Winograd, 1965) to an ancient theorem, the Chinese Remainder Theorem. We will present only an overview of the theory and rely primarily on examples. The interested reader can refer to Dornhoff and Hohn for details.

Out of necessity, integer addition with a computer is addition modulo n , for n some larger number. Consider the case where n is small, like 64. Then addition involves the addition of six-digit binary numbers. Consider the process of adding 31 and 1. Assume the computer's adder takes as input two bit strings $a = \{a_0, a_1, a_2, a_3, a_4, a_5\}$ and $b = \{b_0, b_1, b_2, b_3, b_4, b_5\}$ and outputs $s = \{s_0, s_1, s_2, s_3, s_4, s_5\}$, the sum of a and b . Then, if $a = 31 = (1, 1, 1, 1, 1, 0)$ and $b = 1 = (1, 0, 0, 0, 0, 0)$, s will be $(0, 0, 0, 0, 0, 1)$, or 32. The output $s_5 = 1$ cannot be determined until all other outputs have been determined. If addition is done with a finite-state machine, as in Example 14.3.5, the time required to obtain s will be six time units, where one time unit is the time it takes to get one output from the machine. In general, the time required to obtain s will be proportional to the number of bits. Theoretically, this time can be decreased, but the explanation would require a long digression and our relative results would not change that much. We will use the rule that the number of time units needed to perform addition modulo n is proportional to $\lceil \log_2 n \rceil$.

Now we will introduce a hypothetical problem that we will use to illustrate the idea of a fast adder. Suppose that we had to add many numbers modulo $27720 = 8 \cdot 9 \cdot 5 \cdot 7 \cdot 11$. By the rule above, since $2^{14} < 27720 < 2^{15}$, each addition would take 15 time units. If the sum is initialized to zero, 1,000 additions would be needed; thus, 15,000 time units would be needed to do the additions. We can improve this time dramatically by applying the Chinese Remainder Theorem.

The Chinese Remainder Theorem (CRT). *Let n_1, n_2, \dots, n_p be integers that have no common factor greater than one between any pair of them; i. e., they are relatively prime. Let $n = n_1 n_2 \cdots n_p$. Define*

$$\theta : \mathbb{Z}_n \rightarrow \mathbb{Z}_{n_1} \times \mathbb{Z}_{n_2} \times \cdots \times \mathbb{Z}_{n_p}$$

by

$$\theta(k) = (k_1, k_2, \dots, k_p)$$

where for $1 \leq i \leq p$, $0 \leq k_i < n_i$ and $k \equiv k_i \pmod{n_i}$. Then θ is an isomorphism from \mathbb{Z}_n into $\mathbb{Z}_{n_1} \times \mathbb{Z}_{n_2} \times \cdots \times \mathbb{Z}_{n_p}$.

This theorem can be stated in several different forms, and its proof can be found in many abstract algebra texts.

Example 15.1.8. As we saw in Chapter 11, \mathbb{Z}_6 is isomorphic to $\mathbb{Z}_2 \times \mathbb{Z}_3$. This is the smallest case to which the CRT can be applied. An isomorphism between \mathbb{Z}_6 and $\mathbb{Z}_2 \times \mathbb{Z}_3$ is

$$\begin{aligned} \theta(0) &= (0, 0) & \theta(3) &= (1, 0) \\ \theta(1) &= (1, 1) & \theta(4) &= (0, 1) \\ \theta(2) &= (0, 2) & \theta(5) &= (1, 2) \end{aligned}$$

Let's consider a somewhat larger case. We start by selecting a modulus that can be factored into a product of relatively prime integers.

$$\begin{aligned} n &= 2^5 3^3 5^2 \\ &21600 \end{aligned}$$

In this case the factors are $2^5 = 32$, $3^3 = 27$, and $5^2 = 25$. They need not be powers of primes, but it is easy to break the factors into this form to assure relatively prime numbers. To add in \mathbb{Z}_n , we need $\lceil \log_2 n \rceil = 15$ time units. Let $G = \mathbb{Z}_{32} \times \mathbb{Z}_{27} \times \mathbb{Z}_{25}$. The CRT gives us an isomorphism between \mathbb{Z}_{21600} and G . The basic idea behind the fast adder, illustrated in Figure 15.1.3, is to make use of this isomorphism.

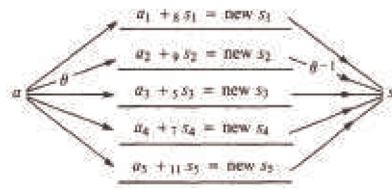


FIGURE 15.1.3

Assume we have several integers a_1, \dots, a_m to be added. Here, we assume $m = 20$.

**a = {1878, 1384, 84, 2021, 784, 1509, 1740, 1201,
2363, 1774, 1865, 33, 1477, 894, 690, 520, 198, 1349, 1278, 650};**

After each of the s_i 's is initialized to zero, each summand t is decomposed into a triple $\theta(t) = (t_1, t_2, t_3) \in G$. For our example we first define θ as a *Mathematica* function and then map it over the list of summands.

$\theta[n_]$:= {Mod[n, 32], Mod[n, 27], Mod[n, 25]}

distributedSummands = Map[θ , a]

**(22 15 3
8 7 9
20 3 9
5 23 21
16 1 9
5 24 9
12 12 15
17 13 1
27 14 13
14 19 24
9 2 15
1 6 8
5 19 2
30 3 19
18 15 15
8 7 20
6 9 23
5 26 24
30 9 3
10 2 0)**

Addition in G can be done in parallel so that each new subtotal in the form of the triple (s_1, s_2, s_3) takes only as long to compute as it takes to add in the largest modulus, $\log_2 32 = 5$ time units, if calculations are done in parallel. By the time rule that we have established, the addition of 20 numbers can be done in $20 \times 5 = 100$ time units, as opposed to $15 \times 20 = 300$ time units if we do the calculations in \mathbb{Z}_n .

The result of adding the distributed summands in the three different moduli for our example would be the following.

**distributedSum =
Fold[{Mod[#1[[1]] + #2[[1]], 32], Mod[#1[[2]] + #2[[2]], 27], Mod[#1[[3]] + #2[[3]], 25]} &, {0, 0, 0}, distributedSummands]
{12, 13, 17}**

Two more factors must still be considered, however. How easy is it to determine $\theta(a)$ and $\theta^{-1}(s_1, s_2, s_3)$? We must compute $\theta(a)$ twenty times, and, if it requires a sizable amount of time, there may not be any advantage to the fast adder. The computation of an inverse is not as time-critical since it must be done only once, after the final sums are determined in G .

The determination of $\theta(a)$ is not a major problem. If the values of $\theta(1)$, $\theta(10)$, $\theta(100)$, $\theta(1000)$, and $\theta(10000)$ are stored, $a = d_0 + 10d_1 + 100d_2 + 1000d_3 + 10000d_4$, then

$$\theta(a) = d_0 \theta(1) + d_1 \theta(10) + d_2 \theta(100) + d_3 \theta(1000) + d_4 \theta(10000)$$

by the fact that θ is an isomorphism. The components of $\theta(a)$ can be computed economically using this formula so as not to slow down the actual adding process.

The computation of $\theta^{-1}(s_1, s_2, s_3)$ is simplified by the fact that θ^{-1} is also an isomorphism. The final sum is $s_1 \theta^{-1}(1, 0, 0) + s_2 \theta^{-1}(0, 1, 0) + s_3 \theta^{-1}(0, 0, 1)$. The arithmetic in this expression is in \mathbb{Z}_{21600} and is more time consuming. However, as was

noted above, it need only be done once. This is why the fast adder is only practical in situations where many additions must be performed to get a single sum.

For our example, we can use Mathematica's built-in function for inverting θ :

```
ChineseRemainder[distributedSum, {32, 27, 25}]
```

```
2092
```

The result we get is exactly what we get by directly adding in the larger modulus.

```
Fold[Mod[#1 + #2, 32 × 27 × 25] &, 0, a]
```

```
2092
```

Notice that if we wanted the conventional sum of integers our list, the result we just arrived at would not be correct. The relationship between the integer sum and the modular sum is that they differ by a multiple of the modulus:

```
Total[a]
```

```
23 692
```

```
Mod[Total[a] - Fold[Mod[#1 + #2, 32 × 27 × 25] &, 0, a], 32 × 27 × 25]
```

```
0
```

To further illustrate the potential of fast adders, consider the problem of addition modulo

$$n = 2^5 \cdot 3^3 \cdot 5^2 \cdot 7^2 \cdot 11 \cdot 13 \cdot 17 \cdot 19 \cdot 23 \cdot 29 \cdot 31 \cdot 37 \cdot 41 \cdot 43 \cdot 47 \approx 3.1 \times 10^{21}$$

Each addition using the usual modulo n addition with full adders would take 72 time units. By decomposing each summand into 15-tuples according to the CRT, the time is reduced to $\lceil \log_2 49 \rceil = 6$ time units per addition.

EXERCISES FOR SECTION 15.1

A Exercises

- What generators besides 1 does $[\mathbb{Z}, +]$ have?
- Without doing any multiplications, determine the number of generators of $[\mathbb{Z}_{11}, +_{11}]$.
- Prove that if $|G| > 2$ and G is cyclic, G has at least two generators.
- If you wanted to list the generators of \mathbb{Z}_n you would only have to test the first $n/2$ positive integers. Why?
- Which of the following groups are cyclic? Explain.
 - $[\mathbb{Q}, +]$
 - $[\mathbb{R}^+, \cdot]$
 - $[6\mathbb{Z}, +]$ where $6\mathbb{Z} = \{6n \mid n \in \mathbb{Z}\}$
 - $\mathbb{Z} \times \mathbb{Z}$
 - $\mathbb{Z}_2 \times \mathbb{Z}_3 \times \mathbb{Z}_{25}$
- For each group and element, determine the order of the cyclic subgroup generated by the element:
 - $\mathbb{Z}_{25}, 15$
 - $\mathbb{Z}_4 \times \mathbb{Z}_9, (2, 6)$ (apply Exercise 8)
 - $\mathbb{Z}_{64}, 2$

B Exercises

- How can Theorem 15.1.4 be applied to list the generators of \mathbb{Z}_n ? What are the generators of \mathbb{Z}_{25} ? Of \mathbb{Z}_{256} ?
- Prove that if the greatest common divisor of n and m is 1, then $(1, 1)$ is a generator of $\mathbb{Z}_n \times \mathbb{Z}_m$, and, hence, $\mathbb{Z}_n \times \mathbb{Z}_m$ is isomorphic to \mathbb{Z}_{nm} .
- Illustrate how the fast adder can be used to add the numbers 21, 5, 7, and 15 using the isomorphism between \mathbb{Z}_{77} and $\mathbb{Z}_7 \times \mathbb{Z}_{11}$.
 - If the same isomorphism is used to add the numbers 25, 26, and 40, what would the result be, why would it be incorrect, and how would the answer differ from the answer in part a?
- Prove that if G is a cyclic group of order n with generator a , and $p, q \in \{0, 1, \dots, n-1\}$, then

$$(p + q)a = (p +_n q)a$$

15.2 Cosets and Factor Groups

Consider the group $[\mathbb{Z}_{12}, +_{12}]$. As we saw in the previous section, we can picture its cyclic properties with the string art of Figure 15.1.1. Here we will be interested in the nongenerators, like 3. The solid lines in Figure 15.2.1 show that only one-third of the tacks have been reached by starting at zero and jumping to every third tack. The numbers of these tacks correspond to $(3) = \{0, 3, 6, 9\}$.

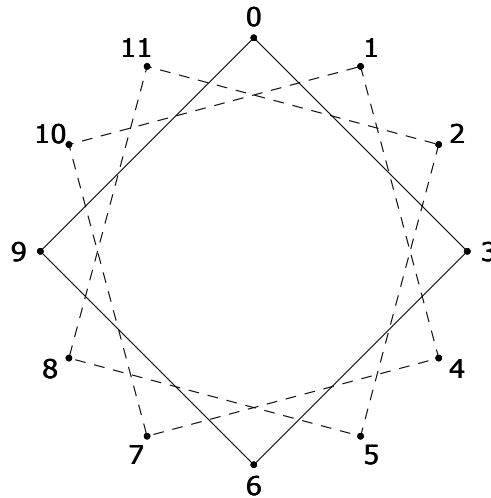


Figure 15.2.1

What happens if you start at one of the unused tacks and again jump to every third tack? The two broken paths on Figure 15.2.1 show that identical squares are produced. The tacks are thus partitioned into very similar subsets. The subsets of \mathbb{Z}_{12} that they correspond to are $\{0, 3, 6, 9\}$, $\{1, 4, 7, 10\}$, and $\{2, 5, 8, 11\}$. These subsets are called *cosets*. In particular, they are called cosets of the subgroup $\{0, 3, 6, 9\}$. We will see that under certain conditions, cosets of a subgroup can form a group of their own. Before pursuing this example any further we will examine the general situation.

Definition: Coset. If $[G, *]$ is a group, $H \leq G$ and $a \in G$, the left coset of H generated by a is

$$a * H = \{a * h \mid h \in H\}.$$

The right coset of H generated by a is

$$H * a = \{h * a \mid h \in H\}.$$

Notes:

- (a) H itself is both a left and right coset since $e * H = H * e = H$.
- (b) If G is abelian, $a * H = H * a$ and the left-right distinction for cosets can be dropped. We will normally use left coset notation in that situation.

Definition: Coset Representative. Any element of a coset is called a representative of that coset.

One might wonder whether a is in any way a special representative of $a * H$ since it seems to define the coset. It is not, as we shall see.

Theorem 15.2.1. If $b \in a * H$, then $a * H = b * H$, and if $b \in H * a$, then $H * a = H * b$.

Remark: A Duality Principle. A duality principle can be formulated concerning cosets because left and right cosets are defined in such similar ways. Any theorem about left and right cosets will yield a second theorem when "left" and "right" are exchanged for "right" and "left."

Proof of Theorem 15.2.1: In light of the remark above, we need only prove the first part of this theorem. Suppose that $x \in a * H$. We need only find a way of expressing x as " b times an element of H ." Then we will have proven that $a * H \subseteq b * H$. By the definition of $a * H$, since b and x are in $a * H$, there exist h_1 and h_2 in H such that $b = a * h_1$ and $x = a * h_2$. Given these two equations, $a = b h_1^{-1}$ and

$$x = a * h_2 = (b * h_1^{-1}) * h_2 = b * (h_1^{-1} * h_2).$$

Since $h_1, h_2 \in H$, $h_1^{-1} * h_2 \in H$, and we are done with this part of the proof. In order to show that $b * H \subseteq a * H$, one can follow essentially the same steps, which we will let the reader fill in. ■

Example 15.2.1. In Figure 15.2.1, you can start at either 1 or 7 and obtain the same path by taking jumps of three tacks in each step. Thus,

$$1 +_{12} \{0, 3, 6, 9\} = 7 +_{12} \{0, 3, 6, 9\} = \{1, 4, 7, 10\}.$$

The set of left (or right) cosets of a subgroup partition a group in a special way:

Theorem 15.2.2. *If $[G, *]$ is a group and $H \leq G$, the set of left cosets of H is a partition of G . In addition, all of the left cosets of H have the same cardinality. The same is true for right cosets.*

Proof: That every element of G belongs to a left coset is clear because $a \in a * H$ for all $a \in G$. If $a * H$ and $b * H$ are left cosets, they are either equal or disjoint. If two left cosets $a * H$ and $b * H$ are not disjoint, $a * H \cap b * H$ is nonempty and some element c belongs to the intersection. Then by Theorem 15.2.1,

$$c \in a * H \Rightarrow a * H = c * H \text{ and}$$

$$c \in b * H \Rightarrow b * H = c * H.$$

Hence $a * H = b * H$.

We complete the proof by showing that each left coset has the same cardinality as H . To do this, we simply observe that if $a \in G$, $\rho : H \rightarrow a * H$ defined by $\rho(h) = a * h$ is a bijection and hence $|H| = |a * H|$. We will leave the proof of this statement to the reader. ■

The function ρ has a nice interpretation in terms of our opening example. If $a \in \mathbb{Z}_n$, the graph of $\{0, 3, 6, 9\}$ is rotated 30° to coincide with one of the three cosets of $\{0, 3, 6, 9\}$.

A Counting Formula. *If $|G| < \infty$ and $H \leq G$, the number of distinct left cosets of H equals $\frac{|G|}{|H|}$. For this reason we use G/H to denote the set of left cosets of H in G .*

Example 15.2.2. The set of integer multiples of four, $4\mathbb{Z}$, is a subgroup of $[\mathbb{Z}, +]$. Four distinct cosets of $4\mathbb{Z}$ partition the integers. They are $4\mathbb{Z}$, $1 + 4\mathbb{Z}$, $2 + 4\mathbb{Z}$, and $3 + 4\mathbb{Z}$, where, for example, $1 + 4\mathbb{Z} = \{1 + 4k \mid k \in \mathbb{Z}\}$. $4\mathbb{Z}$ can also be written $0 + 4\mathbb{Z}$.

Distinguished Representatives. Although we have seen that any representative can describe a coset, it is often convenient to select a distinguished representative from each coset. The advantage to doing this is that there is a unique name for each coset in terms of its distinguished representative. In numeric examples such as the one above, the distinguished representative is usually the smallest nonnegative representative. Remember, this is purely a convenience and there is absolutely nothing wrong in writing $-203 + 4\mathbb{Z}$, $5 + 4\mathbb{Z}$, or $621 + 4\mathbb{Z}$ in place of $1 + 4\mathbb{Z}$ because $-203, 5, 621 \in 1 + 4\mathbb{Z}$.

Before completing the main thrust of this section, we will make note of a significant implication of Theorem 15.2.2. Since a finite group is divided into cosets of a common size by any subgroup, we can conclude:

Lagrange's Theorem. *The order of a subgroup of a finite group must divide the order of the group.*

One immediate implication of Lagrange's Theorem is that if p is prime, \mathbb{Z}_p has no proper subgroups.

We will now describe the operation on cosets which will, under certain circumstances, result in a group. For most of this section, we will assume that G is an abelian group. This is one condition that guarantees that the set of left cosets will form a group.

Definition: Operation on Cosets. *Let C and D be left cosets of H , a subgroup of G with representatives c and d , respectively. Then*

$$C \otimes D = c * H \otimes d * H = (c * d) * H$$

The operation \otimes is called the operation induced on left cosets by $$.*

In Theorem 15.2.3, later in this section, we prove that if G is an abelian group, \otimes is indeed an operation. In practice, if the group G is an additive group, the symbol \otimes is replaced by $+$, as in the following example.

Example 15.2.3. Consider the cosets described in Example 15.2.2. For brevity, we rename $0 + 4\mathbb{Z}$, $1 + 4\mathbb{Z}$, $2 + 4\mathbb{Z}$, and $3 + 4\mathbb{Z}$ with the symbols $\bar{0}$, $\bar{1}$, $\bar{2}$, and $\bar{3}$. Let's do a typical calculation, $\bar{1} + \bar{3}$. We will see that the result is always going to be $\bar{0}$, no matter what representatives we select. For example, $9 \in \bar{1}$, $7 \in \bar{3}$, and $9 + 7 = 16 \in \bar{0}$. Our choice of the representatives $\bar{1}$ and $\bar{3}$ were completely arbitrary. If you are reading this as a *Mathematica* Notebook, you can try out this demonstration that lets you select representatives of the two cosets by moving the sliders and see how the result is consistent.

k1
 k2

Your selection of a representative of $\bar{1}$: 9	Good Choice!
Your selection of a representative of $\bar{3}$: 7	Good Choice!
The sum of representatives is 16	The sum is in $\bar{0}$

Since $C \otimes D$ (or $\bar{1} + \bar{3}$ in this case) can be computed in many ways, it is necessary to show that the choice of representatives does not affect the result. When the result we get for $C \otimes D$ is always independent of our choice of representatives, we say that " \otimes is well defined." Addition of cosets is a well-defined operation on the left cosets of $4\mathbb{Z}$ and is summarized in Table 15.2.1. Do you notice anything familiar?

TABLE 15.2.1
Coset Operation—Table of Example 15.2.3

+	$\bar{0}$	$\bar{1}$	$\bar{2}$	$\bar{3}$
$\bar{0}$	$\bar{0}$	$\bar{1}$	$\bar{2}$	$\bar{3}$
$\bar{1}$	$\bar{1}$	$\bar{2}$	$\bar{3}$	$\bar{0}$
$\bar{2}$	$\bar{2}$	$\bar{3}$	$\bar{0}$	$\bar{1}$
$\bar{3}$	$\bar{3}$	$\bar{0}$	$\bar{1}$	$\bar{2}$

Example 15.2.4. Consider the real numbers, $[\mathbb{R}; +]$, and its subgroup of integers, \mathbb{Z} . Every element of \mathbb{R}/\mathbb{Z} has the same cardinality as \mathbb{Z} . Let $s, t \in \mathbb{R}$. $s \in t + \mathbb{Z}$ if s can be written $t + n$ for some $n \in \mathbb{Z}$. Hence s and t belong to the same coset if they differ by an integer. (See Exercise 6 for a generalization of this fact.)

Now consider the coset $0.25 + \mathbb{Z}$. Real numbers that differ by an integer from 0.25 are 1.25, 2.25, 3.25, ... and -0.75, -1.75, -2.75, If any real number is selected, there exists a representative of its coset that is greater than or equal to 0 and less than 1. We will call that representative the distinguished representative of the coset. For example, 43.125 belongs to the coset represented by 0.125; $-6.382 + \mathbb{Z}$ has 0.618 as its distinguished representative. The operation on \mathbb{R}/\mathbb{Z} is commonly called addition modulo 1. A few typical calculations in \mathbb{R}/\mathbb{Z} are

$$(0.1 + \mathbb{Z}) + (0.48 + \mathbb{Z}) = 0.58 + \mathbb{Z},$$

$$(0.7 + \mathbb{Z}) + (0.31 + \mathbb{Z}) = 0.01 + \mathbb{Z},$$

$$\text{and } -(0.41 + \mathbb{Z}) = -0.41 + \mathbb{Z} = 0.59 + \mathbb{Z}.$$

In general, $-(a + \mathbb{Z}) = (1 - a) + \mathbb{Z}$.

Example 15.2.5. Consider $F = (\mathbb{Z}_4 \times \mathbb{Z}_2)/H$, where $H = \{(0, 0), (0, 1)\}$. Since $\mathbb{Z}_4 \times \mathbb{Z}_2$ is of order 8, each element of F is a coset containing two ordered pairs. We will leave it to the reader to verify that the four distinct cosets are

$$(0, 0) + H, (1, 0) + H, (2, 0) + H, \text{ and } (3, 0) + H.$$

The reader can also verify that F is isomorphic to \mathbb{Z}_4 , since F is cyclic. An educated guess should give you a generator.

Example 15.2.6. Consider the group $\mathbb{Z}_2^4 = \mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_2 \times \mathbb{Z}_2$. Let H be $\langle (1, 0, 1, 0) \rangle$, the cyclic subgroup of \mathbb{Z}_2^4 generated by $(1, 0, 1, 0)$. Since

$$(1, 0, 1, 0) + (1, 0, 1, 0) = (1 +_2 1, 0 +_2 0, 1 +_2 1, 0 +_2 0) = (0, 0, 0, 0)$$

The order of H is 2 and \mathbb{Z}_2^4/H has $|\mathbb{Z}_2^4/H| = \frac{|\mathbb{Z}_2^4|}{|H|} = \frac{16}{2} = 8$ elements. A typical coset is

$$C = (0, 1, 1, 1) + H = \{(0, 1, 1, 1), (1, 1, 0, 1)\}.$$

Since $2(0, 1, 1, 1) = (0, 0, 0, 0)$, $2C = H$, the identity for the operation \mathbb{Z}_2^4/H . The orders of all nonidentity elements of \mathbb{Z}_2^4/H are all 2, and it can be shown that the factor group is isomorphic to \mathbb{Z}_2^3 .

Theorem 15.2.3. If G is an abelian group, and $H \leq G$, the operation induced on cosets of H by the operation of G is well defined.

Proof: Suppose that a, b , and a', b' are two choices for representatives of cosets C and D . That is to say that $a, a' \in C$, $b, b' \in D$. We will show that $a * b$ and $a' * b'$ are representatives of the same coset. Theorem 15.2.1 implies that $C = a * H$ and $D = b * H$, thus we have

$$a' \in a * H \text{ and } b' \in b * H.$$

Then there exists $h_1, h_2 \in H$ such that $a' = a * h_1$ and $b' = b * h_2$ and so

$$\begin{aligned} a' * b' &= (a * h_1) * (b * h_2) \\ &= (a * b) * (h_1 * h_2) \end{aligned}$$

by various group properties and the assumption that H is abelian, which lets us reverse the order in which b and h_1 appear. This last expression for $a' * b'$ implies that $a' * b' \in (a * b) * H$ since $h_1 * h_2 \in H$ because H is a subgroup of G . ■

Theorem 15.2.4. Let G be a group and $H \leq G$. If the operation induced on left cosets of H by the operation of G is well defined, then the set of left cosets forms a group under that operation.

Proof: Let C_1, C_2 , and C_3 be the left cosets with representatives r_1, r_2 , and r_3 , respectively. The values of $C_1 \otimes (C_2 \otimes C_3)$ and $(C_1 \otimes C_2) \otimes C_3$ are determined by $r_1 * (r_2 * r_3)$ and $(r_1 * r_2) * r_3$. By the associativity of $*$ in G , these two group elements are equal and so the two coset expressions must be equal. Therefore, the induced operation is associative. As for the identity and inverse properties, there is no surprise. The identity coset is H , or $e * H$, the coset that contains G 's identity. If C is a coset with representative a , that is, if, $C = a * H$, then C^{-1} is $a^{-1} * H$.

$$(a * H) \otimes (a^{-1} * H) = (a * a^{-1}) * H = e * H = \text{identity coset}.$$

Definition: Factor Group. Let G be a group and $H \leq G$. If the set of left cosets of H forms a group, then that group is called the factor group of G modulo H . It is denoted G/H .

Note: If G is abelian, then every subgroup of G yields a factor group. We will delay further consideration of the non-abelian case to Section 15.4.

Remark on Notation: It is customary to use the same symbol for the operation of G/H as for the operation on G . The reason we used distinct symbols in this section was to make the distinction clear between the two operations.

EXERCISES FOR SECTION 15.2

A Exercises

1. Consider \mathbb{Z}_{10} and the subsets of \mathbb{Z}_{10} , $\{0, 1, 2, 3, 4\}$ and $\{5, 6, 7, 8, 9\}$. Why is the operation induced on these subsets by modulo 10 addition not well defined?
2. Can you think of a group G , with a subgroup H such that $|H| = 6$ and $|G/H| = 6$? Is your answer unique?
3. For each group and subgroup, what is G/H isomorphic to?
 - (a) $G = \mathbb{Z}_4 \times \mathbb{Z}_2$ and $H = \langle (2, 0) \rangle$. Compare to Example 15.2.5.
 - (b) $G = [\mathbb{C}, +]$ and $H = \mathbb{R}$.
 - (c) $G = \mathbb{Z}_{20}$ and $H = \langle 8 \rangle$.
4. For each group and subgroup, what is G/H isomorphic to?
 - (a) $G = \mathbb{Z} \times \mathbb{Z}$ and $H = \{(a, a) \mid a \in \mathbb{Z}\}$.
 - (b) $G = [\mathbb{R}^+, \cdot]$ and $H = \{1, -1\}$.
 - (c) $G = \mathbb{Z}_2^5$ and $H = \langle (1, 1, 1, 1, 1) \rangle$.

B Exercises

5. Prove that if G is a group, $H \leq G$ and $a, b \in G$, $a * H = b * H$ if and only if $b^{-1} * a \in H$.
6. (a) Real addition modulo r , $r > 0$, can be described as the operation induced on cosets of $\langle r \rangle$ by ordinary addition. Describe a system of distinguished representatives for the elements of $\mathbb{R}/\langle r \rangle$.
 (b) Consider the trigonometric function sine. Given that $\sin(x + 2\pi k) = \sin x$ for all $x \in \mathbb{R}$ and $k \in \mathbb{Z}$, show how the distinguished representatives of $\mathbb{R}/\langle 2\pi \rangle$ can be useful in developing an algorithm for calculating the sine of a number.

15.3 Permutation Groups

At the risk of boggling the reader's mind, we will now examine groups whose elements are functions. Recall that a permutation on a set A is a bijection from A into A . Suppose that $A = \{1, 2, 3\}$. There are $3! = 6$ different permutations on A . We will call the set of all 6 permutations S_3 . They are listed in Table 15.3.1. The matrix form for describing a function on a finite set is to list the domain across the top row and the image of each element directly below it. For example $r_1(1) = 2$.

$$\begin{array}{lcl}
 i = \begin{pmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \end{pmatrix} & f_1 = \begin{pmatrix} 1 & 2 & 3 \\ 1 & 3 & 2 \end{pmatrix} \\
 r_1 = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 3 & 1 \end{pmatrix} & f_2 = \begin{pmatrix} 1 & 2 & 3 \\ 3 & 2 & 1 \end{pmatrix} \\
 r_2 = \begin{pmatrix} 1 & 2 & 3 \\ 3 & 1 & 2 \end{pmatrix} & f_3 = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 1 & 3 \end{pmatrix}
 \end{array}$$

Table 15.3.1
Elements of S_3

The operation that will give $\{i, r_1, r_2, f_1, f_2, f_3\}$ a group structure is function composition. Consider the "product" $r_1 \circ f_3$:

$$\begin{aligned}
 r_1 \circ f_3(1) &= r_1(f_3(1)) = r_1(2) = 3 \\
 r_1 \circ f_3(2) &= r_1(f_3(2)) = r_1(1) = 2 \\
 r_1 \circ f_3(3) &= r_1(f_3(3)) = r_1(3) = 1
 \end{aligned}$$

The images of 1, 2, and 3 under $r_1 \circ f_3$ and f_2 are identical. Thus, by the definition of equality for functions, we can say $r_1 \circ f_3 = f_2$. The complete table for the operation of function composition is given in Table 15.3.2. We don't even need the table to verify that we have a group:

- (a) Function composition is always associative (see Chapter 7).
- (b) The identity for the group is i . If g is any one of the permutations on A and $x \in A$,

$$g \circ i(x) = g(i(x)) = g(x)$$

and

$$i \circ g(x) = i(g(x)) = g(x).$$

Therefore $g \circ i = i \circ g = g$.

- (c) A permutation, by definition, is a bijection. In Chapter 7 we proved that this implies that it must have an inverse and the inverse itself is a

bijection and hence a permutation. Hence all elements of S_3 have an inverse in S_3 . If a permutation is displayed in matrix form, its inverse can be obtained by exchanging the two rows and rearranging the columns so that the top row is in order. The first step is actually sufficient to obtain the inverse, but the sorting of the top row makes it easier to recognize the inverse.

Example 15.3.1. Lets consider a typical permutation on $\{1, 2, 3, 4, 5\}$,

$$f = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 5 & 3 & 2 & 1 & 4 \end{pmatrix}.$$

$$f^{-1} = \begin{pmatrix} 5 & 3 & 2 & 1 & 4 \\ 1 & 2 & 3 & 4 & 5 \end{pmatrix} = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 \\ 4 & 3 & 2 & 5 & 1 \end{pmatrix}$$

Note from Table 15.3.2 that this group is non-abelian. Remember, non-abelian is the negation of abelian. The existence of two elements that don't commute is sufficient to make a group non-abelian. In this group, r_1 and f_3 is one such pair: $r_1 \circ f_3 = f_2$ while $f_3 \circ r_1 = f_1$, so $r_1 \circ f_3 \neq f_3 \circ r_1$. Caution: Don't take this to mean that every pair of elements has to have this property. There are several pairs of elements in S_3 that *do* commute. In fact, the identity, i , must commute with everything. Also every element must commute with its inverse.

\circ	i	r_1	r_2	f_1	f_2	f_3
i	i	r_1	r_2	f_1	f_2	f_3
r_1	r_1	r_2	i	f_3	f_1	f_2
r_2	r_2	i	r_1	f_2	f_3	f_1
f_1	f_1	f_2	f_3	i	r_1	r_2
f_2	f_2	f_3	f_1	r_2	i	r_1
f_3	f_3	f_1	f_2	r_1	r_2	i

Table 15.3.2
Operation Table for S_3

Definition: Symmetric Group. Let A be a nonempty set. The set of all permutations on A with the operation of function composition is called the symmetric group on A , denoted S_A .

Our main interest will be in the case where A is finite. The size of A is more significant than the elements, and we will denote by S_k the symmetric group on any set of cardinality k , $k \geq 1$.

Example 15.3.2. Our opening example, S_3 , is the smallest non-abelian group. For that reason, all of its proper subgroups are abelian: in fact, they are all cyclic. Figure 15.3.1 shows the Hasse diagram for the subgroups of S_3 .

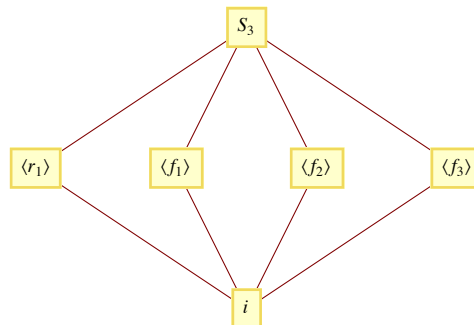


Figure 15.3.1
Lattice diagram of subgroups of S_3

Example 15.3.3. The only abelian symmetric groups are S_1 and S_2 , with 1 and 2 elements, respectively. The elements of S_2 are

$$i = \begin{pmatrix} 1 & 2 \\ 1 & 2 \end{pmatrix} \text{ and } \alpha = \begin{pmatrix} 1 & 2 \\ 2 & 1 \end{pmatrix}$$

S_2 is isomorphic to \mathbb{Z}_2 .

Theorem 15.3.1. For $k \geq 1$, $|S_k| = k!$ and for $k \geq 3$, S_k is non-abelian.

Proof: The first part of the theorem follows from the extended rule of products (see Chapter 2). We leave the details of proof of the second part to the reader after the following hint. Consider f in S_k where $f(1) = 2$, $f(2) = 3$, $f(3) = 1$, and $f(j) = j$ for $3 < j \leq n$. Now define g in a similar manner so that when you compare $f(g(1))$ and $g(f(1))$ you get different results. ■

Cycle Notation

A second way of describing a permutation is by means of cycles, which we will introduce first with an example. Consider $f \in S_8$:

$$f = \begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 8 & 2 & 7 & 6 & 5 & 4 & 1 & 3 \end{pmatrix}$$

Consider the images of 1 when f is applied repeatedly. The images $f(1), f(f(1)), f(f(f(1))), \dots$ are $8, 3, 7, 1, 8, 3, 7, \dots$. If $j \geq 1$, In Figure 15.3.2(a), this situation is represented by the component of the graph that consists of 1, 8, 3, and 7 and shows that the values that you get by repeatedly applying f cycle through those values. This is why we refer to this part of f as a *cycle of length 4*. Of course starting at 8, 3, or 7 also produces the same cycle with only the starting value changing.

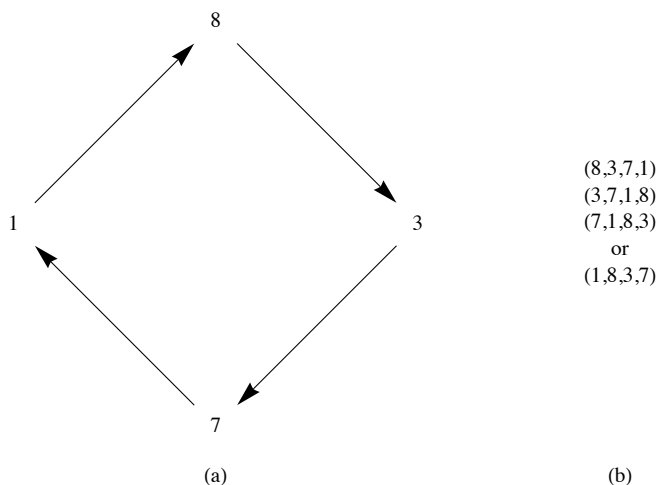


Figure 15.3.2
Representations of cycles of length 4.

Figure 15.3.2(a) illustrates how the cycle can be represented in a visual manner, but it is a bit awkward to write.. Part (b) of the figure presents a more universally recognized way to write a cycle. In (b), a cycle is represented by a list where the image of any number in the list is its successor. In addition, the last number in the list has as its image the first number.

The other elements of the domain of f , are never reached if you start in the cycle $(1, 8, 3, 7)$, and so looking at image of these other numbers will produce numbers that are disjoint from the set $\{1, 8, 3, 7\}$. The other *disjoint cycles* of f are (2) , $(4, 6)$, and (5) . We can express f as a *product of disjoint cycles*:

$$f = (1, 8, 3, 7)(2)(4, 6)(5)$$

or

$$f = (1, 8, 3, 7)(4, 6)$$

where the absence of 2 and 5 implies that $f(2) = 2$ and $f(5) = 5$.

Disjoint Cycles. We say that two cycles are disjoint if no number appears in both cycles, as is the case in our expressions for f above. Disjoint cycles can be written in any order. Thus, we could also say that

$$f = (4, 6)(1, 8, 3, 7).$$

Composing Permutations. We will now consider the composition of permutations written in cyclic form, again by an example. Suppose that $f = (1, 8, 3, 7)(4, 6)$ and $g = (1, 5, 6)(8, 3, 7, 4)$ are elements of S_8 . To calculate $f \circ g$, we start with simple concatenation:

$$f \circ g = (1, 8, 3, 7)(4, 6)(1, 5, 6)(8, 3, 7, 4). \quad (\text{P})$$

Although this is a valid expression for $f \circ g$, our goal is to express the composition as a product of disjoint cycles as f and g were individually written. We will start by determining the cycle that contains 1. *When combining any number of cycles, they are always read from right to left.* The first cycle in (P) does not contain 1; thus we move on to the second. The image of 1 under that cycle is 5. Now we move on to the next cycle, looking for 5, which doesn't appear. The fourth cycle does not contain a 5 either; so $f \circ g(1) = 5$. At this point, we would have written

$$f \circ g = (1, 5$$

on paper. We repeat the steps to determine $f \circ g(5)$. This time the second cycle of (P) moves 5 to 6 and then the third cycle moves 6 to 4. Therefore, $f \circ g(5) = 4$. We continue until the cycle $(1, 5, 4, 3)$ is completed by determining that $f \circ g(3) = 1$. The process is then repeated starting with any number that does not appear in the cycle(s) that have already obtained. The final result for our example is

$$f \circ g = (1, 5, 4, 3)(6, 8, 7).$$

Since $f(2) = 2$ and $g(2) = 2$, $f \circ g(2) = 2$ and we need not include the one-cycle (2) .

Video: For a video that illustrates this process, go to <http://faculty.uml.edu/klevasseur/ads2/videos/cyclecomposition/>.

Example 15.3.4.

$$(a) (1, 2, 3, 4)(1, 2, 3, 4) = (1, 3)(2, 4).$$

$$(b) (1, 4)(1, 3)(1, 2) = (1, 2, 3, 4).$$

Note that the cyclic notation does not indicate the set which is being permuted. The examples above could be in S_5 , where the image of 5 is 5. This ambiguity is usually overcome by making the context clear at the start of a discussion.

Definition: Transposition. A transposition is a cycle of length 2,

Example 15.3.5. $f = (1, 4)$ and $g = (4, 5)$ are transpositions in S_5 . $f \circ g = (1, 4, 5)$ and $g \circ f = (1, 5, 4)$ are not transpositions; thus, the set of transpositions is not closed under composition. Since $f^2 = f \circ f$ and $g^2 = g \circ g$ are both equal to the identity permutation, f and g are their own inverses. In fact, every transposition is its own inverse.

Theorem 15.3.2. Every cycle of length greater than 2 can be expressed as a product of transpositions.

Instead of a formal proof, we will indicate how the product of transpositions can be obtained. The key fact needed is that if $(a_1, a_2, a_3, \dots, a_k)$ is a cycle of length k , it is equal to the following product of $k - 1$ cycles.

$$(a_1, a_k) \cdots (a_1, a_3)(a_1, a_2)$$

Example 11.3.4 (b) illustrates this fact. Of course, a product of cycles can be written as a product of transpositions just as easily by applying the rule above to each cycle. For example,

$$(1, 3, 5, 7)(2, 4, 6) = (1, 7)(1, 5)(1, 3)(2, 6)(2, 4).$$

Unlike the situation with disjoint cycles, we are not free to change the order of these transpositions.

The proofs of the following two theorems appear in many abstract algebra texts.

Theorem 15.3.3. Every permutation on a finite set can be expressed as the product of an even number of transpositions or an odd number of transpositions, but not both.

Theorem 15.3.3 suggests that S_n can be partitioned into its "even" and "odd" elements.

Example 15.3.6. The even permutations of S_3 are i, r_1 and r_2 . They form a subgroup, $\{i, r_1, r_2\}$ of S_3 .

In general:

Theorem 15.3.4. Let $n \geq 2$. The set of even permutations in S_n is a proper subgroup of S_n called the alternating group on $\{1, 2, \dots, n\}$, denoted A_n . The order of A_n is $\frac{n!}{2}$.

Proof: In this proof, the letters s and t stand for transpositions and p, q are even nonnegative integers.

If $f, g \in A_n$, we can write the two permutations as products of even numbers of transpositions:

$$f \circ g = s_1 s_2 \cdots s_p t_1 t_2 \cdots t_q$$

Since $p + q$ is even, $f \circ g \in A_n$. Since A_n is closed With respect to function composition, we have proven that A_n is a subgroup of S_n . by Theorem 11.5.2. To prove the final assertion, let B_n be the set of odd permutations and let $\tau = (1, 2)$. Define $\theta: A_n \rightarrow B_n$ by $\theta(f) = f \circ \tau$. Suppose that $\theta(f) = \theta(g)$. Then $f \circ \tau = g \circ \tau$ and by the cancellation law, $f = g$. Hence, θ is an injection. Next we show that θ is also a surjection. If $h \in B_n$, h is the image of an element of A_n . Specifically, h is the image of $h \circ \tau$.

$$\begin{aligned} \theta(h \circ \tau) &= (h \circ \tau) \circ \tau && \text{Why?} \\ &= h \circ (\tau \circ \tau) && \text{Why?} \\ &= h \circ i && \text{Why?} \\ &= h && \text{Why?} \end{aligned}$$

Since θ is a bijection, $|A_n| = |B_n| = \frac{n!}{2}$. ■

Example 15.3.8. Consider the sliding-tile puzzles pictured in Figure 15.3.3. Each numbered square is a tile and the dark square is a gap. Any tile that is adjacent to the gap can slide into the gap. In most versions of this puzzle, the tiles are locked into a frame so that they can be moved only in the manner described above. The object of the puzzle is to arrange the tiles as they appear in Configuration a. Configurations b and c are typical starting points. We propose to show why the puzzle can be solved starting with b, but not with c.

1	2	3	4	5	6	7	8	5	6	7	8
5	6	7	8	3	4	1	2	3	4	15	2
9	10	11	12	10	9	14	11	10	9	14	11
13	14	15		12	13	15		12	13	1	
(a)				(b)				(c)			

Figure 15.3.3
Configurations of the tile puzzle.

We will associate any configuration of the puzzle with an element of S_{16} . Imagine that a tile numbered 16 fills in the gap. If f is any configuration of the puzzle, i is Configuration a, and for $1 \leq k \leq 16$,

$$f(k) = \text{the number that appears in the position of } k \text{ of } i.$$

If we call Configurations b and c by the names f_1 and f_2 respectively,

$$f_1 = (1, 5, 3, 7)(2, 6, 4, 8)(9, 10)(11, 14, 13, 12)(15)(16)$$

and

$$f_2 = (1, 5, 3, 7, 15)(2, 6, 4, 8)(9, 10)(11, 14, 13, 12)(16).$$

How can we interpret the movement of one tile as a permutation? Consider what happens when the 12 tile of i slides into the gap. The result is a configuration that we would interpret as $(12, 16)$, a single transposition. Now if we slide the 8 tile into the 12 position, the result is $(8, 16, 12)$. Hence, by "exchanging" the tiles 8 and 16, we have obtained $(8, 16)(12, 16) = (8, 16, 12)$.

1	2	3	4
5	6	7	
9	10	11	8
13	14	15	12

Figure 15.3.4
The configuration $(8, 16, 12)$.

Every time you slide a tile into the gap, the new permutation is a transposition composed with the old permutation. Now observe that to start with i and terminate after a finite number of moves with the gap in its original position, you must make an even number of moves. Thus, any permutation that leaves 16 fixed, such as f_1 or f_2 , cannot be solved if it is odd. Note that f_2 is an odd permutation; thus, Puzzle c can't be solved. The proof that all even permutations, such as f_1 , can be solved is left to the interested reader to pursue.

Realizations of Groups. By now we've seen several instances a group can appear through an isomorphic copy of itself in various settings. The simplest such example is the cyclic group of order 2. When this group is mentioned, we might naturally think of the group $[\mathbb{Z}_2, +_2]$, but the groups $[\{-1, 1\}, \cdot]$ and $[S_2, \circ]$ are isomorphic to it. None of these groups are necessarily more natural or important than the others. Which one you use depends on the situation you are in and all are referred to as *realizations* of the cyclic group of order 2. The next family of groups we will study has two natural realizations, first as permutations and second as geometric symmetries.

Example 15.3.9. Dihedral Groups. The dihedral groups can be realized in several ways and we will concentrate on two of them. They can be visualized as symmetries of a regular polygon — this is probably the easiest way to understand the groups. In order to represent the groups on a computer, it is convenient to represent the groups as subgroups of the symmetric groups. If $k \geq 3$, the dihedral group, D_k , is a subgroup of S_k . It is the subgroup of S_k generated by the k -

Realization as symmetries of regular polygons.

We can describe D_n in terms of symmetries of a regular n -gon ($n = 3$: equilateral triangle, $n = 4$: square, $n = 5$: a regular pentagon, ...). Here we will only concentrate on the case of D_4 . If a square is fixed in space, there are several motions of the square that will, at the end of the motion, not change the apparent position of the square. The actual changes in position can be seen if the corners of the square are labeled. In Figure 15.3.5, the initial labeling scheme is shown, along with the four axes of symmetry of the square.

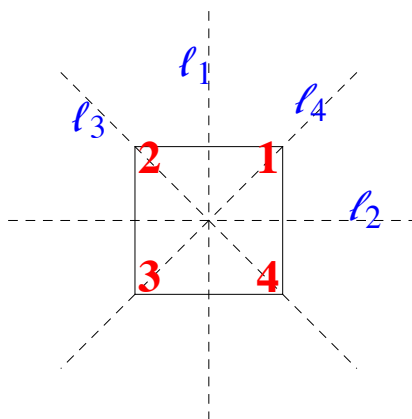


Figure 15.3.5
Axes of symmetry of the square.

It might be worthwhile making a square like this with a sheet of paper. Be careful to label the back so that the numbers match up. Two motions of the square will be considered equivalent if the square is in the same position after performing either motion. There are eight distinct motions. The first four are 0° , 90° , 180° , and 270° clockwise rotations of the square, and the other four are the 180° flips along the axes l_1 , l_2 , l_3 , and l_4 . We

will call the rotations i , r_1 , r_2 , and r_3 , respectively, and the flips f_1 , f_2 , f_3 , and f_4 , respectively. Figure 15.3.6 illustrates r_1 and f_1 . For future

reference we also include the permutations to which they will correspond.

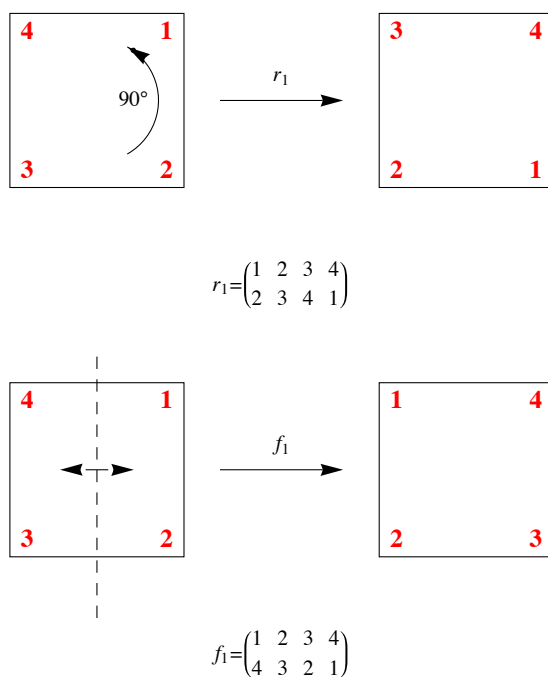


Figure 15.3.6
Two elements of D_4

What is the operation on this set of symmetries? We will call the operation “followed by” and use the symbol $*$ to represent it. The operation will be combine motions, apply motions from right to left, as with functions. We will illustrate how $*$ is computed by finding $r_1 * f_1$. Starting with the initial configuration, if you perform the f_1 motion, and then immediately perform r_1 on the result, we get the same configuration as if we just performed f_4 , which is to flip the square along the line l_4 . Therefore, $r_1 * f_1 = f_4$.

Realization as permutations.

We can also realize the dihedral groups as permutations. For any symmetric motion of the square we can associate with it a permutation. In the case of D_4 , the images of each of the numbers 1 through 4 are the positions on the square that each of the corners 1 through 4 are moved to. For example, since corner 4 moves to position 1 when you perform r_1 , the corresponding function will map 4 to 1. In addition, 1 gets mapped to 2, 2 to 3 and 3 to 4. Therefore, r_1 is the cycle $(1, 2, 3, 4)$. The flip f_1 transposed two pairs of corners and corresponds to $(1, 4)(2, 3)$. If we want to combine these two permutations, using the same names as with motions, we get

$$r_1 \circ f_1 = (1, 2, 3, 4) \circ (1, 4)(2, 3) = (1)(2, 4)(3) = (2, 4).$$

Notice that this permutation corresponds with the flip f_4 .

Although D_4 isn't cyclic (since it isn't abelian), it can be generated from the two elements r_1 and f_1 :

$$D_4 = \langle r_1, f_1 \rangle = \{i, r_1, r_1^2, r_1^3, f_1, r_1 \circ f_1, r_1^2 \circ f_1, r_1^3 \circ f_1\}$$

It is quite easy to describe any of the dihedral groups in a similar fashion. Let

$$r = (1, 2, \dots, n), \text{ an } n\text{-cycle, and}$$

$$f = (1, n)(2, n-1) \dots$$

$$\text{Then } D_n = \langle r, f \rangle = \{i, r, r^2, \dots, r^{n-1}, f, r \circ f, r^2 \circ f, \dots, r^{n-1} \circ f\}$$

An application of D_4 . One application of D_4 is in the design of a letter-facing machine. Imagine letters entering a conveyor belt to be postmarked. They are placed on the conveyor belt at random so that two sides are parallel to the belt. Suppose that a postmarker can recognize a stamp in the top right corner of the envelope, on the side facing up. In Figure 15.3.7, a sequence of machines is shown that will recognize a stamp on any letter, no matter what position in which the letter starts. The letter P stands for a postmarker. The letters R and F stand for rotating and flipping machines that perform the motions of r_1 and f_1 .

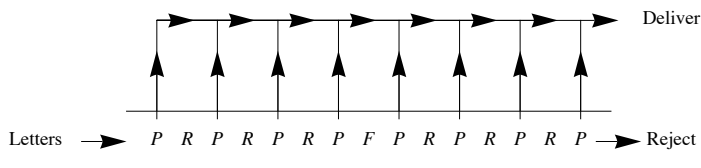


Figure 15.3.7
A letter facer

The arrows pointing up indicate that if a letter is postmarked, it is taken off the conveyor belt for delivery. If a letter reaches the end, it must not have a stamp. Letter-facing machines like this have been designed (see Gallian's paper). One economic consideration is that R-machines tend to cost more than F-machines. R-machines also tend to damage more letters. Taking these facts into consideration, the reader is invited to design a better letter-facing machine. Assume that R-machines cost \$800 and F-machines cost \$500. Be sure that all corners of incoming letters will be examined as they go down the conveyor belt.

EXERCISES FOR SECTION 15.3

A Exercises

1. Given

$$f = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 2 & 1 & 4 & 3 \end{pmatrix}, g = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 2 & 3 & 4 & 1 \end{pmatrix}, \text{ and } h = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 3 & 2 & 4 & 1 \end{pmatrix},$$

compute

- $f \circ g$
- $g \circ h$
- $(f \circ g) \circ h$
- $f \circ (g \circ h)$
- h^{-1}
- $h^{-1} \circ g \circ h$
- f^{-1}

2. Write f , g , and h from Exercise 1 as products of disjoint cycles and determine whether each is odd or even.

3. Do the left cosets of $A_3 = \{i, r_1, r_2\}$ over S_3 form a group under the induced operation on left cosets of A_3 ? What about the left cosets of $\langle f_1 \rangle$?

4. In its realization as permutations, the dihedral group D_3 is equal to S_3 . Can you give a geometric explanation why? Why isn't D_4 equal to S_4 ?

B Exercises

5. (a) Complete the list of elements of D_4 and write out a table for the group in its realization as symmetries.

(b) List the subgroups of D_4 in a lattice diagram. Are they all cyclic? To what simpler groups are the subgroups of D_4 isomorphic?

6. Design a better letter-facing machine (see Example 15.3.9). How can you verify that a letter-facing machine does indeed check every corner of a letter? Can it be done on paper without actually sending letters through it?

7. Prove by induction that if $r \geq 1$ and each t_i is a transposition, then

$$(t_1 \circ t_2 \circ \cdots \circ t_r)^{-1} = t_r \circ \cdots \circ t_2 \circ t_1$$

8. How many elements are there in D_5 ? Describe them geometrically.

9. Complete the proof of Theorem 15.3.1.

10. How many left cosets does A_n , $n \geq 2$ have?

11. Prove that in D_n , $f \circ r = r^{n-1} \circ f$

C Exercise

12. (a) Prove that the tile puzzles corresponding to $A_{16} \cap \{f \in S_{16} \mid f(16) = 16\}$ are solvable.

(b) If $f(16) \neq 16$, how can you determine whether f 's puzzle is solvable?

13. (a) Prove that S_3 is isomorphic to R_3 , the group of 3×3 rook matrices (see Section 11.2 exercises).

(b) Prove that for each $n \geq 2$, R_n is isomorphic to S_n .

15.4 Normal Subgroups and Group Homomorphisms

Our goal in this section is to answer an open question and introduce a related concept. The question is: When are left cosets of a subgroup a group under the induced operation? This question is open for non-abelian groups. Now that we have some examples to work with, we can try a few experiments.

NORMAL SUBGROUPS

Example 15.4.1 $A_3 = \{i, r_1, r_2\}$ is a subgroup of S_3 , and its left cosets are A_3 itself and $B_3 = \{f_1, f_2, f_3\}$. Whether $\{A_3, B_3\}$ is a group boils down to determining whether the induced operation is well defined. Consider the operation table for S_3 in Figure 15.4.1.

\circ	i	r_1	r_2	f_1	f_2	f_3
i	i	r_1	r_2	f_1	f_2	f_3
r_1	r_1	r_2	i	f_3	f_1	f_2
r_2	r_2	i	r_1	f_2	f_3	f_1
f_1	f_1	f_2	f_3	i	r_1	r_2
f_2	f_2	f_3	f_1	r_2	i	r_1
f_3	f_3	f_1	f_2	r_1	r_2	i

Figure 15.4.1
Shaded operation table for S_3

We have shaded in all occurrences of the elements of B_3 in gray. We will call these elements the gray elements and the elements of A_3 the white ones.

Now consider the process of computing the coset product $A_3 \circ B_3$. The "product" is obtained by selecting one white element and one gray element. Note that white "times" gray is always gray. Thus, $A_3 \circ B_3$ is well defined. Similarly, the other three possible products are well defined. The table for the factor group S_3/A_3 is

\circ	A_3	B_3
A_3	A_3	B_3
B_3	B_3	A_3

Clearly, S_3/A_3 is isomorphic to \mathbb{Z}_2 . Note that A_3 and B_3 are also the right cosets of A_3 . This is significant.

Example 15.4.2. Now let's try the left cosets of $\langle f_1 \rangle$ in S_3 . There are three of them. Will we get a complicated version of \mathbb{Z}_3 ? The left cosets are

$$C_0 = \langle f_1 \rangle, C_1 = r_1 \langle f_1 \rangle = \{r_1, f_3\}, \text{ and } C_2 = r_2 \langle f_1 \rangle = \{r_2, f_2\}$$

The reader might be expecting something to go wrong eventually, and here it is. To determine $C_1 \circ C_2$ we can choose from four pairs of representatives:

$$r_1 \in C_1, r_2 \in C_2 \longrightarrow r_1 \circ r_2 = i \in C_0$$

$$r_1 \in C_1, f_2 \in C_2 \longrightarrow r_1 \circ f_2 = f \in C_0$$

$$f_3 \in C_1, r_2 \in C_2 \longrightarrow f_3 \circ r_2 = f_2 \in C_2$$

$$f_3 \in C_1, f_2 \in C_2 \longrightarrow f_3 \circ f_2 = r_2 \in C_2$$

This time, we don't get the same coset for each pair of representatives. Therefore, the induced operation is not well defined and no factor group is obtained.

Commentary: This last development changes our course of action. If we had gotten a factor group from $\{C_0, C_1, C_2\}$, we might have hoped to prove that every collection of left cosets forms a group. Now our question is: How can we determine whether we will get a factor group? Of course, this question is equivalent to: When is the induced operation well defined? There was only one step in the proof of Theorem 15.2.3, where we used the fact that G was abelian. We repeat the equations here:

$$a' * b' = (a * h_1) * (b * h_2) = (a * b) * (h_1 * h_2),$$

since G was abelian.

The last step was made possible by the fact that $h_1 * b = b * h_1$. As the proof continued, we used the fact that $h_1 * h_2$ was in H and so $a' * b'$ is $(a * b) * h$ for some h in H . All that we really needed in the "abelian step" was that

$$h_1 * b = b * (\text{something in } H) = b * h_3.$$

Then, since H is closed under G 's operation, $h_3 * h_2$ is an element of H . The consequence of this observation is included in the following theorem, the proof of which can be found in any abstract algebra text.

Theorem 15.4.1. If $H \leq G$, then the operation induced on left cosets of H by the operation of G is well defined if and only if any one of the following conditions is true:

- (a) If $h \in H, a \in G$, then there exists $h' \in H$ such that $h * a = a * h'$.

(b) If $h \in H, a \in G$, then $a^{-1} * h * a \in H$.

(c) Every left coset of H is equal to a right coset of H .

Corollary 15.4.2. If $H \leq G$, then the operation induced on left cosets of H by the operation of G is well defined if either of the following conditions is true.

(a) G is abelian.

(b) $|H| = \frac{|G|}{2}$.

Example 15.4.3. The right cosets of $\langle f_1 \rangle \leq S_3$ are $\{i, f_1\}$, $\{r_1 f_2\}$, and $\{r_2, f_3\}$. These are not the same as the left cosets of $\langle f_1 \rangle$. In addition, $f_2^{-1} f_1 f_2 = f_2 f_1 f_2 = f_3 \notin \langle f_1 \rangle$.

Definition: Normal Subgroup. If G is a group, $H \leq G$, then H is called a normal subgroup of G , denoted $H \triangleleft G$, if it satisfies any of the conditions of Theorem 15.4.1.

Example 15.4.4. The improper subgroups $\{e\}$ and G of any group G are normal subgroups. $G/\{e\}$ is isomorphic to G . All other normal subgroups of a group, if they exist are called *proper normal subgroups*.

Example 15.4.5. By Condition b of Corollary 15.4.2, A_n is a normal subgroup of S_n and S_n/A_n is isomorphic to \mathbb{Z}_2 .

Example 15.4.6. A_5 , a group in its own right with 60 elements, has many proper subgroups, but none are normal. Although this could be done by brute force, the number of elements in the group would make the process tedious. A far more elegant way to approach the verification of this statement is to use the following fact about the cycle structure of permutations. If $f \in S_n$ is a permutation with a certain cycle structure, $\sigma_1 \sigma_2 \cdots \sigma_k$, where the length of σ_i is ℓ_i , then for any $g \in S_n$, $g^{-1} \circ f \circ g$, which is the conjugate of f by g , will have a cycle structure with exactly the same cycle lengths. For example if we take $f = (1, 2, 3, 4)(5, 6)(7, 8, 9) \in S_9$ and conjugate by $g = (1, 3, 5, 7, 9)$,

$$\begin{aligned} g^{-1} \circ f \circ g &= (1, 9, 7, 5, 3) \circ (1, 2, 3, 4)(5, 6)(7, 8, 9) \circ (1, 3, 5, 7, 9) \\ &= (1, 4, 9, 2)(3, 6)(5, 8, 7) \end{aligned}$$

Notice that the condition for normality of a subgroup H of G is that the conjugate of any element of H by an element of G must remain in H .

To verify that A_5 has no proper normal subgroups, you can start by cataloging the different cycle structures that occur in A_5 and how many elements have those structures. Then consider what happens when you conjugate these different cycle structures with elements of A_5 . An outline of the process is in the exercises.

Example 15.4.7. Let G be the set of two by two invertible matrices of real numbers. That is,

$$G = \left\{ \begin{pmatrix} a & b \\ c & d \end{pmatrix} \mid a, b, c, d \in \mathbb{R}, ad - bc \neq 0 \right\}$$

We saw in Chapter 11 that G is a group with matrix multiplication.

$$H_1 = \left\{ \begin{pmatrix} a & 0 \\ 0 & a \end{pmatrix} \mid a \neq 0 \right\} \text{ and } H_2 = \left\{ \begin{pmatrix} a & 0 \\ 0 & d \end{pmatrix} \mid a, d \neq 0 \right\}$$

are both subgroups of G . H_1 a normal subgroup of G , while H_2 is not normal.

Homomorphisms

Think of the word *isomorphism*. Chances are, one of the first images that comes to mind is an equation something like

$$\theta(x * y) = \theta(x) \diamond \theta(y) \quad (\text{H})$$

An isomorphism must be a bijection, but equation (H) is the algebraic feature of an isomorphism. Here we will examine functions that satisfy equations of this type.

Many homomorphisms are useful since they point out similarities between the two groups (or, on the universal level, two algebraic systems) involved.

Consider the groups $[\mathbb{R}^3, +]$ and $[\mathbb{R}^2, +]$. Every time you use a camera, you are trying to transfer the essence of something three-dimensional onto a photograph—that is, something two-dimensional. If you show a friend a photo you have taken, that person can appreciate much of what you saw, even though a dimension is lacking. The "picture-taking" map is a function $f : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ defined by $f(x_1, x_2, x_3) = (x_1, x_2)$. This function is not a bijection, but it does satisfy the equation $f(x + y) = f(x) + f(y)$ for $x = (x_1, x_2, x_3)$ and $y = (y_1, y_2, y_3)$. Such a function is called a homomorphism, and when a homomorphism exists between two groups, the groups are called homomorphic that is, they are similar. A question that arises with groups, or other algebraic structures, that we claim are homomorphic, or similar, is: How similar are they? When we say that two groups are isomorphic—that is, identical—the map that we use to prove this is unimportant. However, when we say that two groups are homomorphic, the map used gives us a measure of the group's similarities (or dissimilarities). For example, the maps:

$$f_1 : \mathbb{R}^3 \rightarrow \mathbb{R}^3 \text{ defined by } f_1(x_1, x_2, x_3) = (x_1, x_2, x_3),$$

$$f_2 : \mathbb{R}^3 \rightarrow \mathbb{R}^3 \text{ defined by } f_2(x_1, x_2, x_3) = (x_1, x_2, 0), \text{ and}$$

$$f_3 : \mathbb{R}^3 \rightarrow \mathbb{R}^3 \text{ defined by } f_3(x_1, x_2, x_3) = (0, 0, 0)$$

are all homomorphisms. Think of them all as "picture-taking" maps, or cameras. The first camera gives us a three-dimensional picture, the ideal, actually an isomorphism. The second gives us the usual two-dimensional picture, certainly something quite worthwhile. The third

collapses the whole scene onto a point, a "black dot," which gives no idea of the original structure. Hence, the knowledge that two groups are homomorphic doesn't give complete information about the similarities in the structures of the two groups. For this reason, the term homomorphic is rarely used (unlike isomorphic), and the functions, the homomorphisms, are studied.

Definition: Homomorphism. Let $[G, *]$ and $[G', \diamond]$ be groups. $\theta: G \rightarrow G'$ is a homomorphism if $\theta(x * y) = \theta(x) \diamond \theta(y)$ for all $x, y \in G$.

Example 15.4.8. Define $\alpha: \mathbb{Z}_6 \rightarrow \mathbb{Z}_3$ by $\alpha(n) = n(1)$, where $n \in \mathbb{Z}_6$ and $n(1)$ is the sum of n ones in \mathbb{Z}_3 . Therefore, $\alpha(0) = 0$, $\alpha(1) = 1$, $\alpha(2) = 2$, $\alpha(3) = 1 + 1 + 1 = 0$, $\alpha(4) = 1$, and $\alpha(5) = 2$. If $n, m \in \mathbb{Z}_6$,

$$\begin{aligned}\alpha(n +_6 m) &= (n +_6 m)(1) \\ &= n(1) +_3 m(1) \\ &= \alpha(n) +_3 \alpha(m)\end{aligned}$$

Theorem 15.4.2. A few properties of homomorphisms are that if $\theta: G \rightarrow G'$ is a homomorphism, then:

- (a) $\theta(e) = \theta(\text{identity of } G) = \text{identity of } G' = e'$.
- (b) $\theta(a^{-1}) = \theta(a)^{-1}$ for all $a \in G$.
- (c) If $H \leq G$, then $\theta(H) = \{\theta(h) \mid h \in H\} \leq G'$.

Proof:

(a) Let a be any element of G . Then $\theta(a) \in G'$.

$$\begin{aligned}\theta(a) \diamond e' &= \theta(a) && \text{by the definition of } e' \\ &= \theta(a * e) && \text{by the definition of } e \\ &= \theta(a) \diamond \theta(e) && \text{by the fact that } \theta \text{ is a homomorphism}\end{aligned}$$

By cancellation, $e' = \theta(e)$.

(b) Again, let $a \in G$.

$$e' = \theta(e) = \theta(a * a^{-1}) = \theta(a) \diamond \theta(a^{-1}).$$

Hence, by the uniqueness of inverses, $\theta(a)^{-1} = \theta(a^{-1})$.

(c) Let $b_1, b_2 \in \theta(H)$. Then there exists $a_1, a_2 \in H$ such that $\theta(a_1) = b_1$, $\theta(a_2) = b_2$. Recall that a compact necessary and sufficient condition for $H \leq G$ is that $x * y^{-1} \in H$ for all $x, y \in H$. Now we apply the same fact in G' :

$$\begin{aligned}b_1 \diamond b_2^{-1} &= \theta(a_1) \diamond \theta(a_2)^{-1} \\ &= \theta(a_1) \diamond \theta(a_2^{-1}) \\ &= \theta(a_1 * a_2^{-1}) \in \theta(H)\end{aligned}$$

since $a_1 * a_2^{-1} \in H$, and so we can conclude that $\theta(H) \leq G'$. ■

Corollary. Since a homomorphism need not be a surjection and part (c) of Theorem 15.4.2 is true for the case of $H = G$, the range of θ , $\theta(G)$, is a subgroup of G' .

Example 15.4.9. If we define $\pi: \mathbb{Z} \rightarrow \mathbb{Z}/4\mathbb{Z}$ by $\pi(n) = n + 4\mathbb{Z}$, then π is a homomorphism. The image of the subgroup $4\mathbb{Z}$ is the single coset $0 + 4\mathbb{Z}$, the identity of the factor group. Homomorphisms of this type are called *natural homomorphisms*. The following theorems will verify that π is a homomorphism and also show the connection between homomorphisms and normal subgroups. The reader can find more detail and proofs in most abstract algebra texts.

Theorem 15.4.3. If $H \triangleleft G$, then the function $\pi: G \rightarrow G/H$ defined by $\pi(a) = aH$ is a homomorphism, called the *natural homomorphism*.

Based on Theorem 15.4.3, every normal subgroup gives us a homomorphism.

Definition: Kernel. Let $\theta: G \rightarrow G'$ be a homomorphism, and let e' be the identity of G' . The kernel of θ is the set

$$\ker \theta = \{a \in G \mid \theta(a) = e'\}$$

Theorem 15.4.4. Let $\theta: G \rightarrow G'$ be a homomorphism from G into G' . The kernel of θ is a normal subgroup of G .

Based on Theorem 15.4.4, every homomorphism gives us a normal subgroup.

Theorem 15.4.5 : Fundamental Theorem of Group Homomorphisms. Let $\theta: G \rightarrow G'$ be a homomorphism. Then $\theta(G)$ is isomorphic to $G/\ker \theta$.

Example 15.4.10. Define $\theta: \mathbb{Z} \rightarrow \mathbb{Z}_{10}$ by $\theta(n) =$ the remainder from dividing n by 10. The three previous theorems imply the following:

$$(15.4.3) \quad \pi: \mathbb{Z} \rightarrow \mathbb{Z}/10\mathbb{Z} \text{ defined by } \pi(n) = n + 10\mathbb{Z} \text{ is a homomorphism.}$$

$$(15.4.4) \quad \{n \in \mathbb{Z} \mid \theta(n) = 0\} = \{10n \mid n \in \mathbb{Z}\} = 10\mathbb{Z} \triangleleft \mathbb{Z}.$$

$$(15.4.5) \quad \mathbb{Z}/10\mathbb{Z} \text{ is isomorphic to } \mathbb{Z}_{10}.$$

Example 15.4.11. Let G be the same group of two by two invertible real matrices as in Example 15.4.6. Define $\Phi: G \rightarrow G$ by

$\Phi(A) = \frac{A}{\sqrt{|\det A|}}$. We will let the reader verify that Φ is a homomorphism. The theorems above imply:

$$(15.4.4) \quad \ker \Phi = \{A \mid \Phi(A) = I\} = \left\{ \begin{pmatrix} a & 0 \\ 0 & a \end{pmatrix} \mid a \in \mathbb{R}, a \neq 0 \right\} \triangleleft G. \text{ This verifies our statement in Example 15.4.6. As in that example,}$$

let $\ker \Phi = H_1$.

$$(15.4.5) \quad G/H_1 \text{ is isomorphic to } \{A \in G \mid \det A = \pm 1\}.$$

$$(15.4.3) \quad \pi: G \rightarrow G/H_1 \text{ defined, naturally, by } \pi(A) = A H_1 \text{ is a homomorphism.}$$

For the remainder of this section, we will be examining certain kinds of homomorphisms that will play a part in our major application to homomorphisms, coding theory.

Example 15.4.12. Consider $\Phi: \mathbb{Z}_2^2 \rightarrow \mathbb{Z}_2^3$ defined by $\Phi(a, b) = (a, b, a +_2 b)$. If $(a_1, b_1), (a_2, b_2) \in \mathbb{Z}_2^2$,

$$\begin{aligned} \Phi((a_1, b_1) + (a_2, b_2)) &= \Phi(a_1 +_2 a_2, b_1 +_2 b_2) \\ &= (a_1 +_2 a_2, b_1 +_2 b_2, a_1 +_2 a_2 +_2 b_1 +_2 b_2) \\ &= (a_1, b_1, a_1 +_2 b_1) + (a_2, b_2, a_2 +_2 b_2) \\ &= \Phi(a_1, b_1) + \Phi(a_2, b_2) \end{aligned}$$

Since $\Phi(a, b) = (0, 0, 0)$ implies that $a = 0$ and $b = 0$, the kernel of Φ is $\{(0, 0)\}$. By previous theorems, $\Phi(\mathbb{Z}_2^2) = \{(0, 0, 0), (1, 0, 1), (0, 1, 1), (1, 1, 0)\}$ is isomorphic to \mathbb{Z}_2^2 .

We can generalize the previous example as follows: If $n, m \geq 1$ and A an $m \times n$ matrix of 0's and 1's (elements of \mathbb{Z}_2), then $\Phi: \mathbb{Z}_2^m \rightarrow \mathbb{Z}_2^n$ defined by

$$\Phi(a_1, a_2, \dots, a_m) = (a_1, a_2, \dots, a_m) A$$

is a homomorphism. This is true because matrix multiplication is distributive over addition. The only new idea here is that computation is done in \mathbb{Z}_2 where $1 +_2 1 = 0$. If $a = (a_1, a_2, \dots, a_m)$ and $b = (b_1, b_2, \dots, b_m)$, $(a + b)A = aA + bA$ is true by basic matrix laws. Therefore, $\Phi(a + b) = \Phi(a) + \Phi(b)$.

EXERCISES FOR SECTION 15.4

A Exercises

1. Which of the following functions are homomorphisms? What are the kernels of those functions that are homomorphisms?

(a) $\theta_1: \mathbb{R}^* \rightarrow \mathbb{R}^+$ defined by $\theta_1(a) = |a|$.

(b) $\theta_2: \mathbb{Z}_8 \rightarrow \mathbb{Z}_2$ where $\theta_2(n) = \begin{cases} 0 & \text{if } n \text{ is even} \\ 1 & \text{if } n \text{ is odd} \end{cases}$.

(c) $\theta_3: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$, where $\theta_3(a, b) = a + b$.

(d) $\theta_4: S_4 \rightarrow S_4$ defined by $\theta_4(f) = f \circ f = f^2$.

2. Which of the following functions are homomorphisms? What are the kernels of those functions that are homomorphisms?

(a) $\alpha_1: M_{2 \times 2}(\mathbb{R}) \rightarrow \mathbb{R}$, defined by $\alpha_1(A) = A_{11} A_{22} + A_{12} A_{21}$.

(b) $\alpha_2: (\mathbb{R}^*)^2 \rightarrow \mathbb{R}^*$ defined by $\alpha_2(a, b) = ab$.

(c) $\alpha_3: \{A \in M_{2 \times 2}(\mathbb{R}) \mid \det A \neq 0\} \rightarrow \mathbb{R}^*$, where $\alpha_3(A) = \det A$.

(d) $\alpha_4: S_4 \rightarrow S_4$ defined by $\alpha_4(f) = f^{-1}$.

3. Show that D_4 has one proper normal subgroup, but that $\langle (1, 4)(2, 3) \rangle$ is not normal.

4. Prove that the function Φ in Example 15.4.11 is a homomorphism.

5. Define the two functions $\alpha: \mathbb{Z}_2^3 \rightarrow \mathbb{Z}_2^4$ and $\beta: \mathbb{Z}_2^4 \rightarrow \mathbb{Z}_2$ by

$$\alpha(a_1, a_2, a_3) = (a_1, a_2, a_3, a_1 +_2 a_2 +_2 a_3), \text{ and}$$

$$\beta(b_1, b_2, b_3, b_4) = b_1 + b_2 + b_3 + b_4$$

Describe the function $\beta \circ \alpha$. Is it a homomorphism?

6. Express Φ in Example 15.4.12 in matrix form.

B Exercises

7. Prove that if G is an abelian group, then $q(x) = x^2$ defines a homomorphism from G into G . Is q ever an isomorphism?

8. Prove that if $\theta: G \rightarrow G'$ is a homomorphism, and $H \triangleleft G$, then $\theta(H) \triangleleft \theta(G)$. Is it also true that $\theta(H) \triangleleft G'$?

9. Prove that if $\theta : G \rightarrow G'$ is a homomorphism, and $H' \leq \theta(G)$, then $\theta^{-1}(H') = \{a \in G \mid \theta(a) \in H'\} \leq G$.

C Exercises

10. Following up on Example 11.4.6, prove that A_5 is a simple group; i. e., it has no proper normal subgroups.
- (a) Make a list of the different cycle structures that occur in A_5 and how many elements have those structures.
- (b) Within each set of permutations with different cycle structures, identify which subsets are closed with respect to the conjugation operation. With this you will have a partition of A_5 into *conjugate classes* where for each class C ,
- $$f, g \in C \text{ if and only if } \exists \phi \in A_5 \text{ such that } \phi^{-1} \circ f \circ \phi = g$$
- (c) Use the fact that a normal subgroup of A_5 needs to be a union of conjugate classes and verify that no such union exists.

15.5 Coding Theory—Group Codes

In this section, we will introduce the basic ideas involved in coding theory and consider solutions of a coding problem by means of group codes.

A Transmission Problem. Imagine a situation in which information is being transmitted between two points. The information takes the form of high and low pulses (for example, radio waves or electric currents), which we will label 1 and 0, respectively. As these pulses are sent and received, they are grouped together in *blocks* of fixed length. The length determines how much information can be contained in one block. If the length is r , there are 2^r different values that a block can have. If the information being sent takes the form of text, each block might be a character. In that case, the length of a block may be seven, so that $2^7 = 128$ block values can represent letters (both upper and lower case), digits, punctuation, and so on. Figure 15.5.1 illustrates the problem that can be encountered if information is transmitted between two points. During the transmission of data, noise can alter the signal so that what is received differs from what is sent.

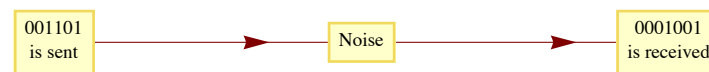


Figure 15.5.1
A noisy transmission

Noise. Noise is a fact of life for anyone who tries to transmit information. Fortunately, in most situations, we could expect a high percentage of the pulses that are sent to be received properly. However, when large numbers of pulses are transmitted, there are usually some errors due to noise. For the remainder of the discussion, we will make assumptions about the nature of the noise and the message that we want to send. Henceforth, we will refer to the pulses as bits.

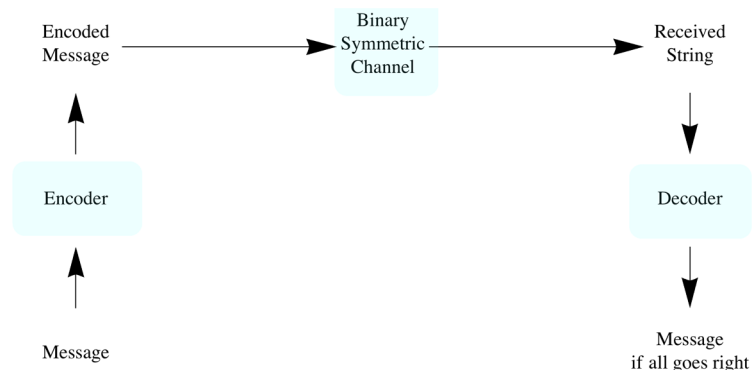


Figure 15.5.2
The Coding Process

Binary Symmetric Channels

We will assume that our information is being sent along a *binary symmetric channel*. By this we mean that any single bit that is transmitted will be received improperly with a certain fixed probability, p . The value of p is usually quite small. To illustrate the process, we will assume that $p = 0.001$, which, in the real world, would be considered somewhat large. Since $1 - p = 0.999$, we can expect 99.9% of all bits to be properly received.

Suppose that our message consists of 3,000 bits of information, to be sent in blocks of three bits each. Two factors will be considered in evaluating a method of transmission. The first is the probability that the message is received with no errors. The second is the number of bits that will be transmitted in order to send the message. This quantity is called the rate of transmission:

$$\text{Rate} = \frac{\text{Message length}}{\text{Number of bits transmitted}}$$

As you might expect, as we devise methods to improve the probability of success, the rate will decrease.

Case 1: Raw information. Suppose that we ignore the noise and transmit the message “as is.” The probability of success is

$$0.999^{3000} = 0.0497124$$

Therefore we only successfully receive the message totally correct less than 5% of the time. The rate of $3000/3000 = 1$ certainly doesn't offset this poor probability.

The Coding Process

Our strategy for improving our chances of success will be to send an encoded message across the binary symmetric channel. The encoding will be done in such a way that small errors can be identified and corrected. This idea is illustrated in Figure 15.5.2.

In our examples, the functions that will correspond to our encoding and decoding devices will all be homomorphisms between Cartesian products of \mathbb{Z}_2 .

Case 2: An Error-Detecting Code. Suppose that each block of three bits $a = (a_1, a_2, a_3)$ is encoded according to the function

$$e : \mathbb{Z}_2^3 \rightarrow \mathbb{Z}_2^4,$$

where

$$e(a) = (a_1, a_2, a_3, a_1 +_2 a_2 +_2 a_3).$$

When the encoded block is received, the first three bits are probably part of the message (it is correct approximately 99.7% of the time), but the added bit that is sent will make it possible to detect single errors in the block. Note that when $e(a)$ is transmitted, the sum of its components is

$$a_1 +_2 a_2 +_2 a_3 +_2 (a_1 +_2 a_2 +_2 a_3) = 0$$

since $a_i + a_i = 0$ in \mathbb{Z}_2 .

If any single bit is garbled by noise, the sum of the received bits will be 1. The last bit of $e(a)$ is called the parity bit. A parity error occurs if the sum of the received bits is 1. Since more than one error is unlikely when p is small, a high percentage of all errors can be detected.

At the receiving end, the decoding function acts on the four-bit block $b = (b_1, b_2, b_3, b_4)$ according to

$$d(b) = (b_1, b_2, b_3, b_1 +_2 b_2 +_2 b_3 +_2 b_4).$$

The fourth bit is called the parity-check bit. If no parity error occurs, the first three bits are recorded as part of the message. If a parity error occurs, we will assume that a retransmission of that block can be requested. This request can take the form of automatically having the parity-check bit of $d(b)$ sent back to the source. If 1 is received, the previous block is retransmitted; if 0 is received, the next block is sent. This assumption of two-way communication is significant, but it is necessary to make this coding system useful. It is reasonable to expect that the probability of a transmission error in the opposite direction is also 0.001. Without going into the details, we will report that the probability of success is approximately 0.990 and the rate is approximately $3/5$. The rate includes the transmission of the parity-check bit to the source.

Case 3: An Error-Correcting Code. For our final case, we will consider a coding process that can correct errors at the receiving end so that only one-way communication is needed. Before we begin, recall that every element of \mathbb{Z}_2^n , $n \geq 1$, is its own inverse; that is, $-b = b$. Therefore, $a - b = a + b$.

The three-bit message blocks are difficult to transmit because they are so similar to one another. If a and b are in \mathbb{Z}_2^3 , their difference, $a +_2 b$, can be thought of as a measure of how close they are. If a and b differ in only one bit position, one error can change one into the other. The encoding that we will introduce takes a block $a = (a_1, a_2, a_3)$ and produces a block of length 6 called the code word of a . The code words are selected so that they are farther from one another than the messages are. In fact, each code word will differ from each other code word by at least three bits. As a result, any single error will not push a code word close enough to another code word to cause confusion. Now for the details. Let

$$G = \begin{pmatrix} 1 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 \end{pmatrix}$$

be the *generator matrix* for the code, and

$$a = (a_1, a_2, a_3)$$

Define $e : \mathbb{Z}_2^3 \rightarrow \mathbb{Z}_2^6$ by

$$e(a) = aG = (a_1, a_2, a_3, a_4, a_5, a_6)$$

where

$$\begin{aligned} a_4 &= a_1 +_2 a_2 \\ a_5 &= a_1 +_2 a_3 \\ a_6 &= a_2 +_2 a_3 \end{aligned}$$

Notice that e is a homomorphism. If a and b are distinct elements of \mathbb{Z}_2^3 , then $c = a + b$ has at least one coordinate equal to 1. Now consider the difference between $e(a)$ and $e(b)$:

$$\begin{aligned} e(a) + e(b) &= e(a + b) \\ &= e(c) \\ &= (c_1, c_2, c_3, c_4, c_5, c_6) \end{aligned}$$

Whether c has 1, 2, or 3 ones, $e(c)$ must have at least three ones; therefore $e(a)$ and $e(b)$ differ in at least three bits.

Now consider the problem of decoding the code words. Imagine that a code word, $e(a)$, is transmitted, and $b = (b_1, b_2, b_3, b_4, b_5, b_6)$ is received. At the receiving end, we know the formula for $e(a)$, and if no error has occurred in transmission,

$$\begin{array}{ll} b_1 = a_1 & \\ b_2 = a_2 & \\ b_3 = a_3 & \\ b_4 = a_1 +_2 a_2 & \Rightarrow \begin{array}{l} b_1 +_2 b_2 +_2 b_4 = 0 \\ b_1 +_2 b_3 +_2 b_5 = 0 \\ b_2 +_2 b_3 +_2 b_6 = 0 \end{array} \\ b_5 = a_1 +_2 a_3 & \\ b_6 = a_2 +_2 a_3 & \end{array}$$

The three equations on the right are called parity-check equations. If any of them is not true, an error has occurred. This error checking can be described in matrix form. Let

$$P = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

P is called the parity-check matrix for this code. Now define $p : \mathbb{Z}_2^6 \rightarrow \mathbb{Z}_2^3$ by $p(b) = bP$. We call $p(b)$ the syndrome of the received block. For example,

$$p(0, 1, 0, 1, 0, 1) = (0, 0, 0) \text{ and } p(1, 1, 1, 1, 0, 0) = (1, 0, 0)$$

Note that p is also a homomorphism. If the syndrome of a block is $(0, 0, 0)$, we can be almost certain that the message block is (b_1, b_2, b_3) .

Next we turn to the method of correcting errors. Despite the fact that there are only eight code words, one for each three-bit block value, the set of possible received blocks is \mathbb{Z}_2^6 , with 64 elements. Suppose that b is not a code word, but that it differs from a code word by exactly one bit. In other words, it is the result of a single error in transmission. Suppose that w is the code word that b is close to and that they differ in the first bit. Then

$$b + w = (1, 0, 0, 0, 0, 0)$$

and

$$\begin{aligned} p(b) &= p(b) + p(w) && \text{since } p(w) = (0, 0, 0) \\ &= p(b + w) && \text{since } p \text{ is a homomorphism} \\ &= p(1, 0, 0, 0, 0, 0) \\ &= (1, 1, 0) \end{aligned}$$

Note that we haven't specified b or w , only that they differ in the first bit. Therefore, if b is received and $p(b) = (1, 1, 0)$, the transmitted code word was probably $b + (1, 0, 0, 0, 0, 0)$ and the message block was $(b_1 +_2 1, b_2, b_3)$. The same analysis can be done if b and w differ in any of the other five bits.

This process can be described in terms of cosets. Let W be the set of code words; that is, $W = e(\mathbb{Z}_2^3)$. W is a subgroup of \mathbb{Z}_2^6 . Consider the factor group \mathbb{Z}_2^6 / W :

$$|\mathbb{Z}_2^6 / W| = \frac{|\mathbb{Z}_2^6|}{|W|} = \frac{64}{8} = 8.$$

Suppose that b_1 and b_2 are representatives of the same coset. Then $b_1 = b_2 + w$ for some w in W . Therefore,

$$\begin{aligned} p(b_1) &= p(b_1) + p(w) && \text{since } p(w) = (0, 0, 0) \\ &= p(b_1 + w) \\ &= p(b_2) \end{aligned}$$

and so b_1 and b_2 have the same syndrome.

Finally, suppose that d_1 and d_2 are distinct and both have only a single coordinate equal to 1. Then $d_1 + d_2$ has exactly two ones. Note that the identity of \mathbb{Z}_2^6 , $(0, 0, 0, 0, 0, 0)$, must be in W . Since $d_1 + d_2$ differs from the identity by two bits, $d_1 + d_2 \notin W$. Hence d_1 and d_2 belong to distinct cosets. The reasoning above serves as a proof of the following theorem.

Theorem 15.5.1. *There is a system of distinguished representatives of \mathbb{Z}_2^6 / W such that each of the six-bit blocks having a single 1 is a distinguished representative of its own coset.*

Now we can describe the error-correcting process. First match each of the blocks with a single 1 with its syndrome. In addition, match the identity of W with the syndrome $(0, 0, 0)$ (see Table 15.5.1). Since there are eight cosets of W , select any representative of the eighth coset to be distinguished. This is the coset with syndrome $(1, 1, 1)$.

Syndrome	Error Correction
0 0 0	0 0 0 0 0 0
1 1 0	1 0 0 0 0 0
1 0 1	0 1 0 0 0 0
0 1 1	0 0 1 0 0 0
1 0 0	0 0 0 1 0 0
0 1 0	0 0 0 0 1 0
0 0 1	0 0 0 0 0 1
1 1 1	1 0 0 0 0 1

Table 15.5.1
Error Correction Table

When block b is received, you need only:

- (1) Compute the syndrome, $p(b)$, and
- (2) Add to b the error correction that matches $p(b)$.

We will conclude this example by computing the probability of success for our hypothetical situation. It is

$$(0.999^6 + 6 \times 0.999^5 \times 0.001)^{1000} = 0.985151.$$

The rate for this method is $\frac{1}{2}$.

EXERCISES FOR SECTION 15.5

A Exercises

1. If the error-detecting code is being used, how would you act on the following received blocks?
 - (a) (1, 0, 1, 1)
 - (b) (1, 1, 1, 1)
 - (c) (0, 0, 0, 0)
2. Express the encoding and decoding functions for the error-detecting code using matrices.
3. If the error-correcting code is being used, how would you decode the following blocks? Expect a problem with one of these. Why?
 - (a) (1, 0, 0, 0, 1, 1)
 - (b) (1, 0, 1, 0, 1, 1)
 - (c) (0, 1, 1, 1, 1, 0)
 - (d) (0, 0, 0, 1, 1, 0)
4. Describe how the triple-repetition code with encoding function, $e: \mathbb{Z}_2 \rightarrow \mathbb{Z}_2^3$, where $e(a_1) = (a_1, a_1, a_1)$ can allow us to correct a single error. What is the probability of success for the $p = 0.001$, 3000-bit situation? What are the generator and parity-check matrices for this code?

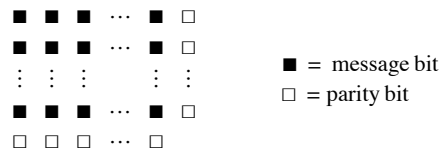
B Exercise

5. Consider the linear code defined the generator matrix

$$G = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{pmatrix}$$

- What size blocks does this code encode and what is the length of the code words?
- What are the code words for this code?
- With this code, can you detect single bit errors? Can you correct all, some, or no single bit errors?

6. **Rectangular codes.** To build a rectangular code, you partition your message into blocks of length m and then factor m into $k_1 \cdot k_2$ and arrange the bits in a k_1 by k_2 rectangular array as in the figure below (read "digit" as "bit"). Then you add parity bits along the right side and bottom of the rows and columns. The code word is read row by row.



For example, if m is 4, then our only choice is a 2 by 2 array. The message 1101 would be encoded as so

1	1	0
0	1	1
1	0	

And the code word is the string 11001110.

- Suppose that you were sent four bit messages using this code and you received the following strings. What were the messages.
 - 11011000
 - 01110010
 - 10001111
- If you encoded n^2 bits in this manner, what would be the rate of the code?
- Rectangular codes are linear codes for the 3 by 2 rectangular code, what are the generator and parity check matrices?

SUPPLEMENTARY EXERCISES FOR CHAPTER 15

Section 15.1

1. How does one find all subgroups of any cyclic group? Can this same process be used to determine all subgroups of noncyclic groups?
2. Exercise 8 of Section 15.1 tells us that $\mathbb{Z}_2 \times \mathbb{Z}_5$ is isomorphic to \mathbb{Z}_{10} . Use the Chinese Remainder Theorem to find an isomorphism between these two groups,
3. Use the Chinese Remainder Theorem to add 74 and 85 in \mathbb{Z}_{120} .

Section 15.2

4. Let G be a group and assume $|G| = 10$. Can G have subgroups of order 2? ...of order 3? ... of order 4? Explain.
5. List all left cosets of $H = \{0, 4, 8\}$ in the group \mathbb{Z}_{12} and write out the table for \mathbb{Z}_{12}/H .
6. Let G be a finite group of order n . Then for any $a \in G$, $a^n = e$, where e is the identity of G . Interpret this statement for the groups $[\mathbb{Z}_6, +_6]$ and $[U(\mathbb{Z}_6), \times_6]$
7. (a) Consider $\mathbb{Z}_8/\langle 2 \rangle$. How many distinct left cosets of $\langle 2 \rangle$ in \mathbb{Z}_8 are there? List them.
(b) Repeat part a for $\mathbb{Z}_{12}/\langle 2 \rangle$.
(c) Is $\mathbb{Z}_8/\langle 2 \rangle$ isomorphic to $\mathbb{Z}_{12}/\langle 2 \rangle$? Explain.

Section 15.3

8. Determine all proper subgroups of the symmetric group S_3 and draw a Hasse diagram for the relation "is a subset of."
9. Let $f \in S_n$. Prove that f is even if and only if f^{-1} is even.
10. (a) By analogy with the motions of a square, how many motions of a cube are there?
(b) Design a "package-facing" machine using the group of motions of the cube.

Section 15.4

11. (a) Let $[B_1, -_1, \vee_1, \wedge_1]$ and $[B_2, -_2, \vee_2, \wedge_2]$ be Boolean algebras. Define a Boolean algebra homomorphism based on the definition of a group homomorphism.
(b) Your definition in part a should result in properties similar to the ones of a group homomorphism. Let $f : B_1 \rightarrow B_2$ be a Boolean algebra homomorphism. Prove:
 - (i) $f(0_1) = 0_2$ and $f(1_1) = 1_2$
 - (ii) $a \leq b \Rightarrow f(a) \leq f(b) \quad \forall a, b \in B_1$ and
 - (iii) $f(B_1)$ is a Boolean subalgebra of B_2 .
12. (a) Prove the contentions of example 15.4.6 that H_1 is a normal subgroup of $GL(2, \mathbb{R})$ but that H_2 is not.
(b) In order to get a clearer picture of what $GL(2, \mathbb{R})/SL(2, \mathbb{R})$ is, prove that the determinant function $\det : GL(2, \mathbb{R}) \rightarrow \mathbb{R}^*$ is an onto homomorphism, and apply Theorem 15.4.5.

Section 15.5

13. This exercise concerns a code called the Hamming (7, 4) code, an error-correcting code with rate $4/7$. A four by seven generator matrix G encodes message blocks of length 4 according to the rule $e(a) = aG$, so that the parity check matrix for the code is

$$P = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

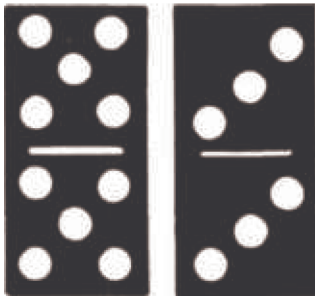
That is, b is a code word iff $bP = (0 \ 0 \ 0)$.

- (a) Find G .
- (b) Encode 1111 and 1001.
- (c) Compute the syndrome of the following received message blocks and correct them, if necessary:
(i) 0100000 (ii) 1010101 (iii) 1011011.
- (d) Prove that this code does indeed correct all single bit errors.

14. Given a code with parity check matrix P whose transpose is given below, identify the generator matrix, and the rate of the code. Prove that the code corrects all single errors.

$$P = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \\ 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

Chapter 16



An Introduction to Rings and Fields

GOALS

In our early elementary school days we began the study of mathematics by learning addition and multiplication on the set of positive integers. We then extended this to operations on the set of all integers. Subtraction and division are defined in terms of addition and multiplication. Later we investigated the set of real numbers under the operations of addition and multiplication. Hence, it is quite natural to investigate those structures on which we can define these two fundamental operations, or operations similar to them. The structures similar to the set of integers are called rings, and those similar to the set of real numbers are called fields.

In coding theory, unstructured coding is at best awkward. Therefore, highly structured codes are needed. The theory of finite fields is essential in the development of such structured codes. We will discuss basic facts about finite fields and introduce the reader to polynomial algebra.

16.1 Rings—Basic Definitions and Concepts

As mentioned in our goals, we would like to investigate algebraic systems whose structure imitates that of the integers.

Definition: Ring. A ring is a set R together with two binary operations, addition and multiplication, denoted by the symbols $+$ and \cdot such that the following axioms are satisfied:

- (1) $[R, +]$ is an abelian group.
- (2) Multiplication is associative on R .
- (3) Multiplication is distributive over addition; that is, for all $a, b, c \in R$, the left distributive law, $a(b + c) = ab + ac$, and the right distributive law, $(b + c)a = ba + ca$, hold.

Comments:

- (1) A ring is designated as $[R, +, \cdot]$ or as just plain R if the operations are understood.
- (2) The symbols $+$ and \cdot stand for arbitrary operations, not just "regular" addition and multiplication. These symbols are referred to by the usual names. For simplicity, we will write ab instead of $a \cdot b$ if it is not ambiguous.
- (3) For the abelian group $[R, +]$, we use additive notation. In particular, the group identity is designated by 0 rather than by e and is customarily called the "zero" of the ring. The group inverse is also written in additive notation: $-a$ rather than a^{-1} .

We now look at some examples of rings. Certainly all the additive abelian groups of Chapter 11 are likely candidates for rings.

Example 16.1.1. $[\mathbb{Z}, +, \cdot]$ is a ring, where $+$ and \cdot stand for regular addition and multiplication on \mathbb{Z} . From Chapter 11, we already know that $[\mathbb{Z}, +]$ is an abelian group, so we need only check parts 2 and 3 of the definition of a ring. From elementary algebra, we know that the associative law under multiplication and the distributive laws are true for \mathbb{Z} . This is our main example of an infinite ring.

Example 16.1.2. $[\mathbb{Z}_n, +_n, \times_n]$ is a ring. The properties of modular arithmetic on \mathbb{Z}_n were described in Section 11.4, and they give us the information we need to convince ourselves that $[\mathbb{Z}_n, +_n, \times_n]$ is a ring. This example is our main example of finite rings of different orders.

Definition: Commutative Ring. A ring in which the commutative law holds under the operation of multiplication is called a commutative ring.

It is common practice to use the word abelian when referring to the commutative law under addition and the word commutative when referring to the commutative law under the operation of multiplication.

Definition: Unity. A ring $[R, +, \cdot]$ that has a multiplicative identity is called a ring with unity. The multiplicative identity itself is called the unity of the ring. More formally, if there exists an element in R , designated by 1, such that for all $x \in R$, $x \cdot 1 = 1 \cdot x = x$, then R is called a ring with unity.

Example 16.1.3. The rings in Examples 16.1.1 and 16.1.2 are commutative rings with unity, the unity in both cases being the number 1.

The ring $[M_{2 \times 2}(\mathbb{R}), +, \cdot]$ is a noncommutative ring with unity, the unity being the identity matrix $I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$.

DIRECT PRODUCTS OF RINGS

Let R_1, R_2, \dots, R_n be rings under the operations $+_1, +_2, \dots, +_n$ and $\cdot_1, \cdot_2, \dots, \cdot_n$ respectively. Let

$$P = \times_{i=1}^n R_i$$

and $a = \{a_1, a_2, \dots, a_n\}, b = \{b_1, b_2, \dots, b_n\} \in P$.

From Chapter 11 we know that P is an abelian group under the operation of componentwise addition:

$$a + b = (a_1 +_1 b_1, a_2 +_2 b_2, \dots, a_n +_n b_n).$$

We also define multiplication on P componentwise:

$$a \cdot b = (a_1 \cdot_1 b_1, a_2 \cdot_2 b_2, \dots, a_n \cdot_n b_n).$$

To show that P is a ring under the above operations, we need only show that the (multiplicative) associative law and the distributive laws hold. This is indeed the case, and we leave it as an exercise. If each of the R_i is commutative, then P is commutative, and if each contains a unity, then P is a ring with unity, which is the n -tuple consisting of the unities of each of the R_i 's.

Example 16.1.4. Since $[\mathbb{Z}_4, +_4, \times_4]$ and $[\mathbb{Z}_3, +_3, \times_3]$ are rings, then $\mathbb{Z}_4 \times \mathbb{Z}_3$ is a ring, where, for example,

$$(2, 1) + (2, 2) = (2 +_4 2, 1 +_3 2) = (0, 0)$$

and

$$(3, 2) \cdot (2, 2) = (3 \times_4 2, 2 \times_3 2) = (2, 1).$$

To determine the unity, if it exists, in the ring $\mathbb{Z}_4 \times \mathbb{Z}_3$, we look for the element (m, n) such that for all elements $(x, y) \in \mathbb{Z}_4 \times \mathbb{Z}_3$,

$$(x, y) = (x, y) \cdot (m, n) = (m, n) \cdot (x, y),$$

or, equivalently,

$$(x \times_4 m, y \times_3 n) = (m \times_4 x, n \times_3 y) = (x, y).$$

So we want m such that $x \times_4 m = m \times_4 x = x$ in the ring \mathbb{Z}_4 . The only element m in \mathbb{Z}_4 that satisfies this equation is $m = 1$. Similarly, we obtain a value of 1 for n . So the unity of $\mathbb{Z}_4 \times \mathbb{Z}_3$, which is unique by Exercise 15 of this section, is $(1, 1)$. We leave to the reader to verify that this ring is commutative.

Hence, products of rings are analogous to products of groups or products of Boolean algebras. We now consider the extremely important concept of multiplicative inverses. Certainly many basic equations in elementary algebra (e.g., $2x = 3$) are solved with this concept. We introduce the main idea here and develop it more completely in the next section.

Example 16.1.5. The equation $2x = 3$ has a solution in the ring $[\mathbb{R}, +, \cdot]$ but does not have a solution in $[\mathbb{Z}, +, \cdot]$, since, to solve this equation, we multiply both sides of the equation $2x = 3$ by the multiplicative inverse of 2. This number, 2^{-1} exists in \mathbb{R} but does not exist in \mathbb{Z} . We formalize this important idea in a definition which by now should be quite familiar to you.

Definition: Multiplicative Inverse. Let $[R, +, \cdot]$ be a ring with unity, 1. If $u \in R$ and there exists an element designated by $v \in R$ such that $u \cdot v = v \cdot u = 1$, then u is said to have a multiplicative inverse, v . We call a ring element that possesses a multiplicative inverse a unit of the ring. The set of all units of a ring R is denoted by $U(R)$.

By Theorem 11.3.2, the multiplicative inverse of a ring element is unique, if it exists. For this reason, we can use the notation u^{-1} for the multiplicative inverse of u , if it exists.

Example 16.1.6. In the rings $[\mathbb{R}, +, \cdot]$ and $[\mathbb{Q}, +, \cdot]$ every nonzero element has a multiplicative inverse. The only elements in \mathbb{Z} that have multiplicative inverses are -1 and 1. That is, $U(\mathbb{R}) = \mathbb{R}^*, U(\mathbb{Q}) = \mathbb{Q}^*$, and $U(\mathbb{Z}) = \{-1, 1\}$.

Example 16.1.7. Let us find the multiplicative inverses, when they exist, of each element of the ring $[\mathbb{Z}_6, +_6, \times_6]$. If $u = 3$, we want an element v such that $u \times_6 v = 1$. We do not have to check whether $v \times_6 u = 1$ since \mathbb{Z}_6 is commutative. If we try each of the six elements, 0, 1, 2, 3, 4, and 5, of \mathbb{Z}_6 , we find that none of them satisfies the above equation, so 3 does not have a multiplicative inverse in \mathbb{Z}_6 . However, since $5 \times_6 5 = 1$, 5 does have a multiplicative inverse in \mathbb{Z}_6 , namely itself: $5^{-1} = 5$. The following table summarizes all results for \mathbb{Z}_6 .

u	u^{-1}
0	does not exist
1	1
2	does not exist
3	does not exist
4	does not exist
5	5

It shouldn't be a surprise that the zero of a ring is never going to have a multiplicative inverse except in the trivial case of $R = \{0\}$.

Isomorphism is a universal concept that is important in every algebraic structure. Two rings are isomorphic as rings if and only if they have the same cardinality and if they behave exactly the same under corresponding operations. They are essentially the same ring. For this to be true, they must behave the same as groups (under $+$) and they must behave the same under the operation of multiplication.

Definition: Ring Isomorphism. Let $[R, +, \cdot]$ and $[R', +', \cdot']$ be rings. Then R is isomorphic to R' if and only if there exists a map, $f: R \rightarrow R'$, called a ring isomorphism, such that

- (1) f is one-to-one and onto,
- (2) $f(a + b) = f(a) + ' f(b)$ for all $a, b \in R$, and
- (3) $f(a \cdot b) = f(a) \cdot ' f(b)$ for all $a, b \in R$.

Conditions 1 and 2 tell us that f is a group isomorphism. Therefore, to show that two rings are isomorphic, we must produce a map, called an isomorphism, that satisfies the definition. Sometimes it is quite difficult to find a map that works. This does not necessarily mean that no such isomorphism exists, but simply that we cannot find it.

This leads us to the problem of how to show that two rings are not isomorphic. This is a universal concept. It is true for any algebraic structure and was discussed in Chapter 11. To show that two rings are not isomorphic, we must demonstrate that they behave differently under one of the operations. We illustrate through several examples.

Example 16.1.8. Consider the rings $[\mathbb{Z}, +, \cdot]$ and $[2\mathbb{Z}, +, \cdot]$. In Chapter 11 we showed that as groups, the two sets \mathbb{Z} and $2\mathbb{Z}$ with addition were isomorphic. The group isomorphism that proved this was the map $f: \mathbb{Z} \rightarrow 2\mathbb{Z}$, defined by $f(n) = 2n$. Is f a ring isomorphism? We need only check whether $f(m \cdot n) = f(m) \cdot ' f(n)$ for all $m, n \in \mathbb{Z}$:

$$f(m \cdot n) = 2 \cdot m \cdot n \text{ and}$$

$$f(m) \cdot ' f(n) = 2m \cdot 2n = 4 \cdot m \cdot n$$

Therefore, f is not a ring isomorphism. This does not necessarily mean that the two rings \mathbb{Z} and $2\mathbb{Z}$ are not isomorphic, but simply that the f doesn't satisfy the conditions. We could imagine that some other function does. We could proceed and try to determine another function f to see whether it is a ring isomorphism, or we could try to show that \mathbb{Z} and $2\mathbb{Z}$ are not isomorphic as rings. To do the latter, we must find something different about the ring structure of \mathbb{Z} and $2\mathbb{Z}$.

We already know that they behave identically under addition, so if they are different as rings, it must have something to do with how they behave under the operation of multiplication. Let's begin to develop a checklist of how the two rings could differ:

- (1) Do they have the same cardinality? Yes, they are both countable.
- (2) Are they both commutative? Yes.
- (3) Are they both rings with unity? No.

\mathbb{Z} is a ring with unity, namely the number 1. $2\mathbb{Z}$ is not a ring with unity, $1 \notin 2\mathbb{Z}$. Hence, they are not isomorphic as rings.

Example 16.1.9. Next consider whether $[2\mathbb{Z}, +, \cdot]$ and $[3\mathbb{Z}, +, \cdot]$ are isomorphic. Because of the previous example, we might guess that they are not. However, checklist items 1 through 3 above do not help us. Why? We add another checklist item:

- (4) Find an equation that makes sense in both rings, which is solvable in one and not the other.

The equation $x + x = x \cdot x$, or $2x = x^2$, makes sense in both rings. However, this equation has a nonzero solution, $x = 2$, in $2\mathbb{Z}$, but does not have a nonzero solution in $3\mathbb{Z}$. Thus we have an equation solvable in one ring that cannot be solved in the other, so they cannot be isomorphic.

Another universal concept that applies to the theory of rings is that of a subsystem. A subring of a ring $[R, +, \cdot]$ is any nonempty subset S of R that is a ring under the operations of R . First, for S to be a subring of the ring R , S must be a subgroup of the group $[R, +]$. Also, S must be closed under \cdot , satisfy the associative law (under \cdot), and satisfy the distributive laws. But since R is a ring, the associative and distributive laws are true for every element in R , and, in particular, for all elements in S , since $S \subseteq R$. We have just proven the following theorem:

Theorem 16.1.1. A subset S of a ring $[R, +, \cdot]$ is a subring of R if and only if:

- (1) $[S, +]$ is a subgroup of the group $[R, +]$, which by Theorem 11.5.1, means we must show:
 - (a) If $a, b \in S$, then $a + b \in S$,
 - (b) $0 \in S$, and
 - (c) If $a \in S$, then $-a \in S$.

(2) S is closed under multiplication: if $a, b \in S$, then $a \cdot b \in S$.

Example 16.1.10. The set of even integers, $2\mathbb{Z}$, is a subring of the ring $[\mathbb{Z}, +, \cdot]$ since $[2\mathbb{Z}, +]$ is a subgroup of the group $[\mathbb{Z}, +]$ and since it is also closed with respect to multiplication:

$$2m, 2n \in 2\mathbb{Z} \Rightarrow (2m) \cdot (2n) = 2(2 \cdot m \cdot n) \in 2\mathbb{Z}.$$

Several of the basic facts that we are familiar with are true for any ring. The following theorem lists a few of the elementary properties of rings.

Theorem 16.1.2. Let $[R, +, \cdot]$ be a ring, with $a, b \in R$. Then

- (1) $a \cdot 0 = 0 \cdot a = 0$
- (2) $a \cdot (-b) = (-a) \cdot b = -(a \cdot b)$
- (3) $(-a) \cdot (-b) = a \cdot b$

Proof of Part 1:

$$\begin{aligned} a \cdot 0 &= a \cdot (0 + 0) \\ &= a \cdot 0 + a \cdot 0 \text{ by the left distributive law.} \end{aligned}$$

Hence if we add $-(a \cdot 0)$ to both sides of the above, we obtain $a \cdot 0 = 0$. Similarly, we can prove that $0 \cdot a = 0$.

Proof of Part 2: Before we begin the proof of part 2, recall that the inverse of each element of the group $[R, +]$ is unique. Hence the inverse of the element $a \cdot b$ is unique and it is denoted $-(a \cdot b)$.

Therefore, to prove that $a \cdot (-b) = -(a \cdot b)$, we need only show that $a \cdot (-b)$ inverts $a \cdot b$.

$$\begin{aligned} a \cdot (-b) + a \cdot b &= a \cdot (-b + b) \text{ by the distributive axiom} \\ &= a \cdot 0 \text{ since } -b \text{ inverts } b \\ &= 0 \text{ by part 1 of this theorem} \end{aligned}$$

Similarly, it can be shown that $(-a) \cdot b = -(a \cdot b)$. This completes the proof of part 2.

We leave the proof of part 3 to the reader (see Exercise 16 of this section). ■

Example 16.1.11. We will compute $2 \cdot (-2)$ in the ring $[\mathbb{Z}_6, +_6, \times_6]$.

$$2 \times_6 (-2) = -(2 \times_6 2) = -4 = 2,$$

since the additive inverse of 4 (mod 6) is 2. Of course, we could have done the calculation directly as

$$2 \times_6 (-2) = 2 \times_6 4 = 2.$$

As the example above illustrates, Theorem 16.1.2 is a modest beginning in the study of which algebraic manipulations are possible in the solution of problems in rings. A fact in elementary algebra that is used frequently in problem solving is the cancellation law. We know that the cancellation laws are true under addition for any ring (Theorem 11.3.5).

Are the cancellation laws true under multiplication? More specifically, let $[R, +, \cdot]$ be a ring and let $a, b, c \in R$ with $a \neq 0$. When can we cancel the a 's in the equation $a \cdot b = a \cdot c$? We can certainly do so if a^{-1} exists, but we cannot assume that a has a multiplicative inverse. The answer to this question is found with the following definition and Theorem 16.1.3.

Definition: Divisors of Zero. Let $[R, +, \cdot]$ be a ring. If a and b are two nonzero elements of R such that $a \cdot b = 0$, then a and b are called divisors of zero.

Example 16.1.12 (a) In the ring $[\mathbb{Z}_8, +_8, \times_8]$, the numbers 4 and 2 are divisors of zero since $4 \times_8 2 = 0$. In addition, 6 is a divisor of zero because $6 \times_8 4 = 0$.

(b) In the ring $[M_{2 \times 2}(\mathbb{R}), +, \cdot]$ the matrices $A = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$ and $B = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$ are divisors of zero since $AB = 0$.

Example 16.1.13. $[\mathbb{Z}, +, \cdot]$ has no divisors of zero.

Now here is why divisors of zero are related to cancellation.

Theorem 16.1.3. The (multiplicative) cancellation law holds in a ring $[R, +, \cdot]$ if and only if R has no divisors of zero.

We prove the theorem using the left cancellation law, namely that if $a \neq 0$ and $a \cdot b = a \cdot c$, then $b = c$ for all $a, b, c \in R$. The proof is similar using the right cancellation law.

Proof: (\Rightarrow) Assume the left cancellation law holds in R and assume that a and b are two elements in R such that $a \cdot b = 0$. We must show that either $a = 0$ or $b = 0$. To do this, assume that $a \neq 0$ and show that b must be 0.

$$\begin{aligned} a \cdot b = 0 &\Rightarrow a \cdot b = a \cdot 0 \text{ by Theorem 16.2.1, part 1} \\ &\Rightarrow b = 0 \text{ by the cancellation law} \end{aligned}$$

(\Leftarrow) Conversely, assume that R has no divisors of 0 and we will prove that the cancellation law must hold. To do this, assume that $a, b, c \in R$, $a \neq 0$, such that $a \cdot b = a \cdot c$ and show that $b = c$.

$$\begin{aligned}
 a \cdot b = a \cdot c &\Rightarrow a \cdot b - a \cdot c = 0 && \text{Why?} \\
 &\Rightarrow a \cdot (b - c) = 0 && \text{Why?} \\
 &\Rightarrow b - c = 0 && \text{Why?} \\
 &\Rightarrow b = c && \blacksquare
 \end{aligned}$$

Hence, the only time that the cancellation laws hold in a ring is when there are no divisors of zero. The commutative rings with unity in which the above is true are given a special name.

Definition: Integral Domain. A commutative ring with unity containing no divisors of zero is called an integral domain.

In this chapter, Integral domains will be denoted generically by the letter D .

We state the following two useful facts without proof.

Theorem 16.1.4. The element m in the ring \mathbb{Z}_n is a divisor of zero if and only if m is not relatively prime to n (i.e., $\gcd(m, n) \neq 1$).

Corollary. If p is a prime, then \mathbb{Z}_p has no divisors of zero.

Example 16.1.14. $[\mathbb{Z}, +, \cdot]$, $[\mathbb{Z}_p, +_p, \times_p]$ with p a prime, $[\mathbb{Q}, +, \cdot]$, $[\mathbb{R}, +, \cdot]$, and $[\mathbb{C}, +, \cdot]$ are all integral domains. The key example of an infinite integral domain is $[\mathbb{Z}, +, \cdot]$. In fact, it is from \mathbb{Z} that the term integral domain is derived. The main example of a finite integral domain is $[\mathbb{Z}_p, +_p, \times_p]$, when p is prime.

We close this section with the verification of an observation that was made in Chapter 11, namely that the product of two algebraic systems may not be an algebraic system of the same type.

Example 16.1.15. Both $[\mathbb{Z}_2, +_2, \times_2]$ and $[\mathbb{Z}_3, +_3, \times_3]$ are integral domains. Consider the product $\mathbb{Z}_2 \times \mathbb{Z}_3$. It's true that $\mathbb{Z}_2 \times \mathbb{Z}_3$ is a commutative ring with unity (see Exercise 13). However, $(1, 0) \cdot (0, 2) = (0, 0)$, so $\mathbb{Z}_2 \times \mathbb{Z}_3$ has divisors of zero and is therefore not an integral domain.

EXERCISES FOR SECTION 16.1

A Exercises

1. Review the definition of rings to show that the following are rings. The operations involved are the usual operations defined on the sets. Which

of these rings are commutative? Which are rings with unity? For the rings with unity, determine the unity and all units.

- (a) $[\mathbb{Z}, +, \cdot]$
- (b) $[\mathbb{C}, +, \cdot]$
- (c) $[M_{n \times n}(\mathbb{R}), +, \cdot]$
- (d) $[\mathbb{Q}, +, \cdot]$
- (e) $[M_{2 \times 2}(\mathbb{R}), +, \cdot]$
- (f) $[\mathbb{Z}_2, +_2, \times_2]$

2. Follow the instructions for Exercise 1 and the following rings:

- (a) $[\mathbb{Z}_6, +_6, \times_6]$
- (b) $[\mathbb{Z}_5, +_5, \times_5]$
- (c) $[\mathbb{Z}_2^3, +, \cdot]$
- (d) $[\mathbb{Z}_8, +_8, \times_8]$
- (e) $[\mathbb{Z} \times \mathbb{Z}, +, \cdot]$
- (f) $[\mathbb{R}^2, +, \cdot]$

3. Show that the following pairs of rings are not isomorphic:

- (a) $[\mathbb{Z}, +, \cdot]$ and $[M_{2 \times 2}(\mathbb{Z}), +, \cdot]$
- (b) $[3\mathbb{Z}, +, \cdot]$ and $[4\mathbb{Z}, +, \cdot]$.

4. Show that the following pairs of rings are not isomorphic:

- (a) $[\mathbb{R}, +, \cdot]$ and $[\mathbb{Q}, +, \cdot]$.
- (b) $[\mathbb{Z}_2 \times \mathbb{Z}_2, +, \cdot]$ and $[\mathbb{Z}_4, +, \cdot]$.

5. (a) Show that $3\mathbb{Z}$ is a subring of the ring $[\mathbb{Z}, +, \cdot]$
- (b) Find all subrings of \mathbb{Z}_8 .
- (c) Find all subrings of $\mathbb{Z}_2 \times \mathbb{Z}_2$.
6. Verify the validity of Theorem 16.1.3 by finding examples of elements a, b , and c ($a \neq 0$) in the following rings, where $a \cdot b = a \cdot c$ and yet $b \neq c$:
- (a) \mathbb{Z}_8
- (b) $M_{2 \times 2}(\mathbb{R})$
- (c) \mathbb{Z}_2^2
7. (a) Determine all solutions of the equation $x^2 - 5x + 6 = 0$ in \mathbb{Z} . Can there be any more than two solutions to this equation (or any quadratic equation) in \mathbb{Z} ?
- (b) Find all solutions of the equation in part a in \mathbb{Z}_{12} . Why are there more than two solutions?
8. Solve the equation $x^2 + 4x + 4 = 0$ in the following rings. Interpret 4 as $1 + 1 + 1 + 1$, where 1 is the unity of the ring.
- (a) in \mathbb{Z}_8
- (b) in $M_{2 \times 2}(\mathbb{R})$
- (c) in \mathbb{Z}
- (d) in \mathbb{Z}_3

B Exercises

9. The relation “is isomorphic to” on rings is an equivalence relation. Explain the meaning of this statement.
10. Let R_1, R_2, \dots, R_n be rings. Prove the multiplicative, associative, and distributive laws for the ring
- $$R = \prod_{i=1}^n R_i$$
- (a) If each of the R_i is commutative, is R commutative?
- (b) Under what conditions will R be a ring with unity?
- (c) What will the units of R be when it has a unity?
11. (a) Prove that the ring $\mathbb{Z}_2 \times \mathbb{Z}_3$ is commutative and has unity.
- (b) Determine all divisors of zero for the ring $\mathbb{Z}_2 \times \mathbb{Z}_3$.
- (c) Give another example illustrating the fact that the product of two integral domains may not be an integral domain. Is there an example where the product is an integral domain?
12. **Boolean Rings.** Let U be a nonempty set.
- (a) Verify that $[\mathcal{P}(U), \oplus, \cap]$ is a commutative ring with unity.
- (b) What are the units of this ring?
13. (a) For any ring $[R, +, \cdot]$, expand $(a + b)(c + d)$ for $a, b, c, d \in R$.
- (b) If R is commutative, prove that $(a + b)^2 = a^2 + 2ab + b^2$ for all $a, b \in R$.
14. (a) Let R be a commutative ring with unity. Prove by induction that for $n \geq 1$,
- $$(a + b)^n = \sum_{k=0}^n \binom{n}{k} a^k b^{n-k}$$
- (b) Simplify $(a + b)^5$ in \mathbb{Z}_5 .
- (c) Simplify $(a + b)^{10}$ in \mathbb{Z}_{10} .
15. Prove: If R is a ring with unity then this unity is unique.
16. Prove part 3 of Theorem 16.1.2.
17. Prove the Corollary to Theorem 16.1.4.

18. Let U be a finite set. Prove that the Boolean ring $[\mathcal{P}(U), \oplus, \cap]$ is isomorphic to the ring $[\mathbb{Z}_2^n, +, \cdot]$, where $n = |U|$

16.2 Fields

Although the algebraic structures of rings and integral domains are widely used and play an important part in the applications of mathematics, we still cannot solve the simple equation $ax = b$, $a \neq 0$ in all rings or in all integral domains. Yet this is one of the first equations we learn to solve in elementary algebra and its solvability is basic to innumerable questions. Certainly, if we wish to solve a wide range of problems in a system we need at least all of the laws true for rings and the cancellation laws together with the ability to solve the equation $ax = b$, $a \neq 0$. We summarize the above in a definition and list several theorems without proof that will place this concept in the context of the previous section.

Definition: Field. A field is a commutative ring with unity such that each nonzero element has a multiplicative inverse.

In this chapter, we denote a field generically by the letter F . The letters k , K and L are also conventionally used for fields.

Example 16.2.1. $[\mathbb{Q}, +, \cdot]$, $[\mathbb{R}, +, \cdot]$, and $[\mathbb{C}, +, \cdot]$ are all fields.

Reminder: Since every field is a ring, all facts and concepts that are true for rings are true for any field.

Theorem 16.2.1. Every field is an integral domain.

Of course the converse of Theorem 16.2.1 is not true. Consider $[\mathbb{Z}, +, \cdot]$.

Theorem 16.2.2. Every finite integral domain is a field.

Theorem 16.2.3. If p is a prime, then \mathbb{Z}_p is a field.

Theorem 16.2.3 is immediate from Theorem 16.2.2.

Theorem 16.2.1 reminds us that the cancellation laws must be true for any field. Theorem 16.2.3 gives us a large number of finite fields, but we must be cautious. This theorem does not tell us that all finite fields are of the form \mathbb{Z}_p , p a prime. To see this, let's try to construct a field of order 4.

Example 16.2.2: a field of order 4. First the field must contain the additive and multiplicative identities, 0 and 1, so, without loss of generality, we can assume that the field we are looking for is of the form $F = \{0, 1, a, b\}$. Since there are only two nonisomorphic groups of order 4, we have only two choices for the group table for $[F, +]$. If the additive group is isomorphic to \mathbb{Z}_4 then two of the nonzero elements of F would not be their own additive inverse (as are 1 and 3 in \mathbb{Z}_4). Let's assume $\beta \in F$ is one of those elements and $\beta + \beta = \gamma \neq 0$. An isomorphism between the additive groups F and \mathbb{Z}_4 would require that γ in F correspond with 2 in \mathbb{Z}_4 . We could continue our argument and infer that $\gamma \cdot \gamma = 0$, producing a zero divisor, which we need to avoid if F is to be a field. We leave the remainder of the argument to the reader. We can thus complete the addition table so that $[F, +]$ is isomorphic to \mathbb{Z}_2^2 :

+	0	1	a	b
0	0	1	a	b
1	1	0	b	a
a	a	b	0	1
b	b	a	1	0

Next, by Theorem 16.1.2, Part 1, and since 1 is the unity of F , the table for multiplication must look like:

·	0	1	a	b
0	0	0	0	0
1	0	1	a	b
a	0	a	—	—
b	0	b	—	—

Hence, to complete the table, we have only four entries to find, and, since F must be commutative, this reduces our task to filling in three entries. Next, each nonzero element of F must have a unique multiplicative inverse. The inverse of a must be either a itself or b . If $a^{-1} = a$, then $b^{-1} = b$. (Why?) But

$a^{-1} = a \Rightarrow a \cdot a = 1$. And if $a \cdot a = 1$, then $a \cdot b$ is equal to a or b . In either case, by the cancellation law, we obtain $a = 1$ or $b = 1$, which is impossible. Therefore we are forced to conclude that $a^{-1} = b$ and $b^{-1} = a$. To determine the final two products of the table, simply note that $a \cdot a \neq a$ because the equation $x^2 = x$ has only two solutions, 0 and 1 in any field. We also know that $a \cdot a$ cannot be 1 because a doesn't invert itself and cannot be 0 because a can't be a zero divisor. This leaves us with one possible conclusion, that $a \cdot a = b$ and similarly $b \cdot b = a$. Hence, our multiplication table for F is:

·	0	1	a	b
0	0	0	0	0
1	0	1	a	b
a	0	a	b	1
b	0	b	1	a

The table listing the multiplicative inverse of each nonzero element is:

u	u^{-1}
1	1
a	b
b	a

We leave it to the reader to convince him- or herself, if it is not already clear, that $[F, +, \cdot]$, as described above, is a field. Hence, we have produced a field of order 4 and 4 is not a prime.

This construction would be difficult to repeat for larger fields. In section 16.4 we will introduce a different approach to constructing fields that will be far more efficient.

Even though not all finite fields are isomorphic to \mathbb{Z}_p , for some prime p it can be shown that every field F must have either:

- (1) a subfield isomorphic to \mathbb{Z}_p for some prime p , or
- (2) a subfield isomorphic to \mathbb{Q} .

In particular, if F is a finite field, a subfield of F must exist that is isomorphic to \mathbb{Z}_p . One can think of all fields as being constructed from either \mathbb{Z}_p or \mathbb{Q} .

Example 16.2.3. $[\mathbb{R}, +, \cdot]$ is a field, and it contains a subfield isomorphic to $[\mathbb{Q}, +, \cdot]$, namely \mathbb{Q} itself.

Example 16.2.4. The field F that we constructed in Example 16.2.2 should have a subfield isomorphic to \mathbb{Z}_p for some prime p . From the tables, we note that the subset $\{0, 1\}$ of $\{0, 1, a, b\}$ under the given operations of F behaves exactly like $[\mathbb{Z}_2, +_2, \times_2]$. Hence, the field in Example 16.2.2 has a subfield isomorphic to \mathbb{Z}_2 . Does it have a subfield isomorphic to a larger field, say \mathbb{Z}_3 ? We claim not and leave this investigation to the reader (see Exercise 3 of this section).

We close this section with a brief discussion of isomorphic fields. Again, since a field is a ring, the definition of isomorphism of fields is the same as that of rings. It can be shown that if f is a field isomorphism, then $f(a^{-1}) = f(a)^{-1}$; that is, inverses are mapped onto inverses under any field isomorphism. A major question to try to solve is: How many different non-isomorphic finite fields are there of any given order? If p is a prime, it seems clear from our discussions that all fields of order p are isomorphic to \mathbb{Z}_p . But how many nonisomorphic fields are there, if any, of order 4, 6, 8, 9, etc? The answer is given in the following theorem, whose proof is beyond the scope of this text.

Theorem 16.2.4.

- (1) Any finite field F has order p^n for a prime p and a positive integer n .
- (2) For any prime p and any positive integer n there is a field of order p^n .
- (3) Any two fields of order p^n are isomorphic. This field of order p^n is frequently referred to as the Galois field of order p^n and it is designated by $GF(p^n)$.

Evariste Galois (1811-32) was a pioneer in the field of abstract algebra.



A French stamp honoring Evariste Galois (1811-32)

Theorem 16.2.4 tells us that there is a field of order $2^2 = 4$, and there is only one such field up to isomorphism. That is, all such fields of order 4 are isomorphic to F , which we constructed in Example 16.2.2.

EXERCISES FOR SECTION 16.2

A Exercises

1. Write out the addition, multiplication, and "inverse" tables for each of the following fields'.

- (a) $[\mathbb{Z}_2, +_2, \times_2]$
- (b) $[\mathbb{Z}_3, +_3, \times_3]$
- (c) $[\mathbb{Z}_5, +_5, \times_5]$
2. Show that the set of units of the fields in Exercise 1 form a group under the operation of the multiplication of the given field. Recall that a unit is an element which has a multiplicative inverse.
3. Complete the argument in Example 16.2.2 to show that if $[F, +]$ is isomorphic to \mathbb{Z}_4 , then F would have a zero divisor.
4. Write out the operation tables for \mathbb{Z}_2^2 . Is \mathbb{Z}_2^2 a ring? An integral domain? A field? Explain.
5. Determine all values x from the given field that satisfy the given equation:
- (a) $x + 1 = -1$ over $\mathbb{Z}_2, \mathbb{Z}_3$ and \mathbb{Z}_5
- (b) $2x + 1 = 2$ over \mathbb{Z}_3 and \mathbb{Z}_5
- (c) $3x + 1 = 2$ over \mathbb{Z}_5
6. (a) Prove that if p and q are prime, then $\mathbb{Z}_p \times \mathbb{Z}_q$ is never a field.
- (b) Can \mathbb{Z}_p^n be a field for any prime p and any positive integer $n \geq 2$?
7. The following are equations over \mathbb{Z}_2 . Their coefficients come solely from \mathbb{Z}_2 . Determine all solutions over \mathbb{Z}_2 ; that is, find all numbers in \mathbb{Z}_2 that satisfy the equations:
- (a) $x^2 + x = 0$
- (b) $x^2 + 1 = 0$
- (c) $x^3 + x^2 + x + 1 = 0$
- (d) $x^3 + x + 1 = 0$
8. Determine the number of different fields, if any, of all orders 2 through 15. Wherever possible, describe these fields via a known field.

B Exercise

9. Let $\mathbb{Q}(\sqrt{2}) = \{a + b\sqrt{2} \mid a, b \in \mathbb{Q}\}$.
- (a) Prove that $[\mathbb{Q}(\sqrt{2}), +, \cdot]$ is a field.
- (b) Show that \mathbb{Q} is a subfield of $\mathbb{Q}(\sqrt{2})$. For this reason, $\mathbb{Q}(\sqrt{2})$ is called an extension field of \mathbb{Q} .
- (c) Show that all the roots of the equation $x^2 - 2 = 0$ lie in the extension field $\mathbb{Q}(\sqrt{2})$.
- (d) Do the roots of the equation $x^2 - 3 = 0$ lie in this field? Explain.

16.3 Polynomial Rings

In the previous sections we examined the solutions of a few equations over different rings and fields. To solve the equation $x^2 - 2 = 0$ over the field of the real numbers means to find all solutions of this equation that are in this particular field \mathbb{R} . This statement can be replaced as follows: Determine all $a \in \mathbb{R}$ such that the polynomial $f(x) = x^2 - 2$ is equal to zero when evaluated at $x = a$. In this section, we will concentrate on the theory of polynomials. We will develop concepts using the general setting of polynomials over rings since results proven over rings are true for fields (and integral domains). The reader should keep in mind that in most cases we are just formalizing concepts that he or she learned in high school algebra over the field of reals.

Definition: Polynomial over R . Let $[R, +, \cdot]$ be a ring. A polynomial, $f(x)$, over R is an expression of the form

$$f(x) = \sum_{i=0}^n a_i x^i = a_0 + a_1 x + a_2 x^2 + \cdots + a_n x^n, \quad n \geq 0,$$

where $a_0, a_1, a_2, \dots, a_n \in R$. If $a_n \neq 0$, then the degree of $f(x)$ is n . If $f(x) = 0$, then the degree of $f(x)$ is undefined and we assign the value $-\infty$ to the degree. If the degree of $f(x)$ is n , we write $\deg f(x) = n$.

Comments:

- (1) The symbol x is an object called an *indeterminate*, which is not an element of the ring R .
- (2) The set of all polynomials in the indeterminate x with coefficients in R is denoted by $R[x]$.
- (3) Note that $R \subseteq R[x]$. The elements of R are called *constant polynomials*, with the nonzero elements of R being the polynomials of degree 0.
- (4) R is called the *ground ring* for $R[x]$.
- (5) In the definition above, we have written the terms in increasing degree starting with the constant. The ordering of terms can be reversed without changing the polynomial. For example, $1 + 2x - 3x^4$ and $-3x^4 + 2x + 1$ are the same polynomial.
- (6) A term of the form x^k in a polynomial is understood to be $1x^k$.

Example 16.3.1. $f(x) = 3$, $g(x) = 2 - 4x + 7x^2$, and $h(x) = 2 + x^4$ are all polynomials in $\mathbb{Z}[x]$. Their degrees are 0, 2, and 4, respectively.

Addition and multiplication of polynomials are performed as in high school algebra. However, we must do our computations in the ground ring over which we are considering the polynomials.

Example 16.3.2. In $\mathbb{Z}_3[x]$, if $f(x) = 1 + x$ and $g(x) = 2 + x$, then

$$\begin{aligned} f(x) + g(x) &= (1 + x) + (2 + x) \\ &= (1 +_3 2) + (1 +_3 1)x \\ &= 0 + 2x \\ &= 2x \end{aligned}$$

and

$$\begin{aligned} f(x)g(x) &= (1 + x) \cdot (2 + x) \\ &= (1 + x) \cdot 2 + (1 + x) \cdot x \\ &= 1 \times_3 2 + 2x + 1x + x \cdot x \\ &= 2 + (2 +_3 1)x + x^2 \\ &= 2 + x^2 \end{aligned}$$

However, for the same polynomials as above, $f(x)$ and $g(x)$ in $\mathbb{Z}[x]$, we have

$$\begin{aligned} f(x) + g(x) &= (1 + x) + (2 + x) \\ &= (1 + 2) + (1 + 1)x \\ &= 3 + 2x \end{aligned}$$

and

$$\begin{aligned} f(x)g(x) &= (1 + x) \cdot (2 + x) \\ &= (1 + x) \cdot 2 + (1 + x) \cdot x \\ &= 1 \cdot 2 + 2x + 1x + x \cdot x \\ &= 2 + (2 + 1)x + x^2 \\ &= 2 + 3x + x^2 \end{aligned}$$

The important fact to keep in mind is that addition and multiplication in $R[x]$ depends on addition and multiplication in R . The x 's merely serve the purpose of "place holders." All computations are done over the given ring. We summarize in the following theorem:

Theorem 16.3.1. Let $[R, +, \cdot]$ be a ring. Then:

- (1) $R[x]$ is a ring under the operations of polynomial addition and multiplication, which depend on (are induced by) the operations in R .
- (2) If R is a commutative ring, then $R[x]$ is a commutative ring.

(3) If R is a ring with unity, 1, then $R[x]$ is a ring with unity (the unity in $R[x]$ is $1 + 0x + 0x^2 + \dots$).

(4) If R is an integral domain, then $R[x]$ is an integral domain.

(5) If F is a field, then $F[x]$ is not a field. However, $F[x]$ is an integral domain.

The proofs for Parts 1 through 4 are not difficult but rather long, so we omit them. For those inclined to prove them, we include the formal definitions of addition and multiplication in $R[x]$ below.

Proof Of Part 5: $F[x]$ is not a field since for $x \in F[x]$, $x^{-1} = 1/x \notin F[x]$. Hence not all nonzero elements in $F[x]$ have multiplicative inverses in $F[x]$. Every field F is an integral domain. By Part 4, $F[x]$ is an integral domain. ■

Definition: Addition in $R[x]$. Let $f(x) = a_0 + a_1x + a_2x^2 + \dots + a_mx^m$ and $g(x) = b_0 + b_1x + b_2x^2 + \dots + b_nx^n$ be elements in $R[x]$ so that $a_i \in R$ and $b_i \in R$ for all i . Let k be the maximum of m and n . Then

$$f(x) + g(x) = c_0 + c_1x + c_2x^2 + \dots + c_kx^k$$

where $c_i = a_i + b_i$ for $i = 0, 1, 2, \dots, k$.

Definition: Multiplication in $R[x]$. Let $f(x)$ and $g(x)$ be as above. Then

$$f(x) \cdot g(x) = d_0 + d_1x + d_2x^2 + \dots + d_px^p \text{ where}$$

$p = m + n$, and

$$\begin{aligned} d_s &= \sum_{i=0}^s a_i b_{s-i} \\ &= a_0 b_s + a_1 b_{s-1} + a_2 b_{s-2} + \dots + a_{s-1} b_1 + a_s b_0 \end{aligned}$$

for $0 \leq s \leq p$.

Example 16.3.3. Let $f(x) = 2 + x^2$ and $g(x) = -1 + 4x + 3x^2$. We will compute $f(x) \cdot g(x)$ in $\mathbb{Z}[x]$. Of course this product can be obtained by the usual methods of high school algebra. We will, for illustrative purposes, use the above definition. Using the notation of the above definition, $a_0 = 2$, $a_1 = 0$, $a_2 = 1$, $b_0 = -1$, $b_1 = 4$, and $b_2 = 3$. We want to compute the coefficients d_0, d_1, d_2, d_3 , and d_4 . We will compute d_3 , the coefficient of the x^3 term of the product, and leave the remainder to the reader (see Exercise 2 of this section). Since the degrees of both factors is 2, $a_i = b_i = 0$ for $i \geq 3$.

$$\begin{aligned} d_3 &= a_0 b_3 + a_1 b_2 + a_2 b_1 + a_3 b_0 \\ &= 2 \cdot 0 + 0 \cdot 3 + 1 \cdot 4 + 0 \cdot (-1) = 4 \end{aligned}$$

From high school algebra we all learned the standard procedure for dividing a polynomial $f(x)$ by a second polynomial $g(x)$. This process of polynomial long division is referred to as the division property for polynomials. Under this scheme we continue to divide until the result is a quotient $q(x)$ and a remainder $r(x)$ whose degree is strictly less than that of the divisor $g(x)$. This property is valid over any field.

Example 16.3.4. Let $f(x) = 1 + x + x^3$ and $g(x) = 1 + x$ be two polynomials in $\mathbb{Z}_2[x]$. Let us divide $f(x)$ by $g(x)$. Keep in mind that we are in $\mathbb{Z}_2[x]$ and that, in particular, $-1 = 1$ in \mathbb{Z}_2 . This is a case where reordering the terms in decreasing degree is preferred.

$$\begin{array}{r} x^2 + x \\ x+1 \overline{) x^3 + 0x^2 + x + 1} \\ \underline{x^3 + x^2} \\ x^2 + x + 1 \\ \underline{x^2 + x} \\ 1 \end{array}$$

Therefore,

$$\frac{x^3 + x + 1}{x+1} = x^2 + x + \frac{1}{x+1}$$

or equivalently,

$$x^3 + x + 2 = (x^2 + x) \cdot (x+1) + 1$$

That is $f(x) = g(x) \cdot q(x) + r(x)$ where $q(x) = x^2 + x$ and $r(x) = 1$. Notice that $\deg(r(x)) = 0$, which is strictly less than the $\deg(g(x)) = 1$.

Example 16.3.5. Let $f(x) = 1 + x^4$ and $g(x) = 1 + x$ be polynomials in $\mathbb{Z}_2[x]$. Let us divide $f(x)$ by $g(x)$:

$$\begin{array}{r}
 x^3 + x^2 + x + 1 \\
 x+1 \overline{) x^4 + 0x^3 + 0x^2 + 0x + 1} \\
 \underline{x^4 + x^3} \\
 x^3 \\
 \underline{x^3 + x^2} \\
 x^2 \\
 \underline{x^2 + x} \\
 x + 1 \\
 \underline{x + 1} \\
 0
 \end{array}$$

Thus $x^4 + 1 = (x^3 + x^2 + x + 1)(x + 1)$.

Since we have 0 as a remainder, $x + 1$ must be a factor of $x^4 + 1$, as in high school algebra. Also, since $x + 1$ is a factor of $x^4 + 1$, 1 is a zero (or root) of $x^4 + 1$. Of course we could have determined that 1 is a root of $f(x)$ simply by computing $f(1) = 1^4 + 1 = 1 + 1 = 0$.

Before we summarize the main results suggested by the previous examples, we should probably consider what could have happened if we had performed divisions of polynomials in the ring $\mathbb{Z}[x]$ rather than over the field \mathbb{Z}_2 . For example, $f(x) = x^2 - 1$ and $g(x) = 2x - 2$ are both elements of the ring $\mathbb{Z}[x]$, yet

$$\frac{x^2+1}{2x-1} = \frac{1}{2}x + \frac{1}{2}$$

The quotient is not a polynomial over \mathbb{Z} but a polynomial over the field \mathbb{Q} . For this reason it would be wise to describe all results over a field F rather than over an arbitrary ring R .

Theorem 16.3.2. Division Property for $F[x]$. Let $[F, +, \cdot]$ be a field and let $f(x)$ and $g(x)$ be two elements of $F[x]$ with $g(x) \neq 0$. Then there exist unique polynomials $q(x)$ and $r(x)$ in $F[x]$ such that $f(x) = g(x)q(x) + r(x)$, where $\deg r(x) < \deg g(x)$.

Theorem 16.3.2 can be proven by induction on $\deg f(x)$.

Theorem 16.3.3. Let $[F, +, \cdot]$ be a field. An element $a \in F$ is a zero of $f(x) \in F[x]$ if and only if $x - a$ is a factor of $f(x)$ in $F[x]$.

Proof: (\Rightarrow) Assume that $a \in F$ is a zero of $f(x) \in F[x]$. We wish to show that $x - a$ is a factor of $f(x)$. To do so, apply the division property to $f(x)$ and $g(x) = x - a$. Hence, there exist unique polynomials $q(x)$ and $r(x)$ from $F[x]$ such that $f(x) = (x - a) \cdot q(x) + r(x)$ and the $\deg r(x) < \deg(x - a) = 1$, so $r(x) = c \in F$, that is, $r(x)$ is a constant. Also a is a zero of $f(x)$ means $f(a) = 0$. So $f(x) = (x - a) \cdot q(x) + c$ becomes $0 = f(a) = (a - a)q(a) + c$. Hence $c = 0$, so $f(x) = (x - a) \cdot q(x)$, and $x - a$ is a factor of $f(x)$. The reader should note that a critical point of the proof of this half of the theorem was the part of the division property that stated that $\deg r(x) < \deg g(x)$.

(\Leftarrow) We leave this half to the reader, exercise 6. ■

Theorem 16.3.4. A nonzero polynomial $f(x) \in F[x]$ of degree n can have at most n zeros.

Proof: Let $a \in F$ be a zero of $f(x)$. Then $f(x) = (x - a) \cdot q(x)$ by Theorem 16.3.3. If $b \in F$ is a zero of $q(x)$, then again by Theorem 16.3.3, $f(x) = (x - a)(x - b)q(x)$. Continue this process, which must terminate in at most n steps. ■

From Theorem 16.3.3 we can obtain yet another insight into the problems associated with solving polynomial equations; that is, finding the zeros of a polynomial. The theorem states that an element $a \in F$ is a zero of $f(x) \in F[x]$ if and only if $x - a$ is a factor of $f(x)$. The initial important idea here is that the zero a is from the ground field F . Second, a is a zero only if $(x - a)$ is a factor of $f(x)$ in $F[x]$ —that is, only when $f(x)$ can be factored (or reduced) to the product of $(x - a)$ times some other polynomial in $F[x]$.

Example 16.3.6. Consider the polynomial $f(x) = x^2 - 2$ taken as being in $\mathbb{Q}[x]$. From high school algebra we know that $f(x)$ has two zeros (or roots), namely $\pm\sqrt{2}$, and $x^2 - 2$ can be factored (reduced) as $(x - \sqrt{2})(x + \sqrt{2})$. However, we are working in $\mathbb{Q}[x]$, these two factors are not in the set of polynomials over the rational numbers, \mathbb{Q} since $\sqrt{2} \notin \mathbb{Q}$. Therefore, $x^2 - 2$ does not have a zero in \mathbb{Q} since it cannot be factored over \mathbb{Q} . When this happens, we say that the polynomial is irreducible over \mathbb{Q} .

The problem of factoring polynomials is tied hand-in-hand with that of the reducibility of polynomials. We give a precise definition of this concept.

Definition: Irreducible over F . Let $[F, +, \cdot]$ be a field and let $f(x) \in F[x]$ be a nonconstant polynomial, $f(x)$ is irreducible over F if and only if $f(x)$ cannot be expressed as a product of two (or more) polynomials, both from $F[x]$ and both of degree lower than that of $f(x)$.

A polynomial is reducible over F if it is not irreducible over F .

Example 16.3.7. The polynomial $f(x) = x^4 + 1$ of Example 16.3.5 is reducible over \mathbb{Z}_2 since $x^4 + 1 = (x + 1)(x^3 + x^2 + x + 1)$.

Example 16.3.8. Is the polynomial $f(x) = x^3 + x + 1$ of Example 16.3.4 reducible over \mathbb{Z}_2 ? From Example 16.3.4 we know that $x + 1$ is not a factor of $x^3 + x + 1$, and from high school algebra we realize that a cubic (also second-degree) polynomial is reducible if and only if it has a linear (first-degree) factor. (Why?) Does $f(x) = x^3 + x + 1$ have any other linear factors? Theorem 16.3.1 gives us a quick way of determining this since $x - a$ is a factor of $x^3 + x + 1$ over \mathbb{Z}_2 if and only if $a \in \mathbb{Z}_2$ is a zero of $x^3 + x + 1$. So $x^3 + x + 1$ is reducible over \mathbb{Z}_2 if and only if it has a zero in \mathbb{Z}_2 . Since \mathbb{Z}_2 has only two elements, 0 and 1, this is easy enough to check.

$$f(0) = 0^3 + 0 + 1 = 1 \quad \text{and}$$

$$f(1) = 1^3 + 1 + 1 = 1$$

so neither 0 nor 1 is a zero of $f(x)$ over \mathbb{Z}_2 . Hence, $x^3 + x + 1$ is irreducible over \mathbb{Z}_2 .

From high school algebra we know that $x^3 + x + 1$ has three zeros from some field. Can we find this field? To be more precise, can we construct (find) the field which contains \mathbb{Z}_2 and all zeros of $x^3 + x + 1$? We will consider this task in the next section.

We close this section with a final analogy. Prime numbers play an important role in mathematics. The concept of irreducible polynomials (over a field) is analogous to that of a prime number. Just think of the definition of a prime number. A useful fact concerning primes is: If p is a prime and if $p \mid a \cdot b$, then $p \mid a$ or $p \mid b$. We leave it to the reader to think about the veracity of the following: If $p(x)$ is an irreducible polynomial over F , $a(x), b(x) \in F[x]$ and $p(x) \mid a(x) \cdot b(x)$, then $p(x) \mid a(x)$ or $p(x) \mid b(x)$.

EXERCISES FOR SECTION 16.3

A Exercises

1. Let $f(x) = 1 + x$ and $g(x) = 1 + x + x^2$. Compute the following sums and products in the indicated rings.

(a) $f(x) + g(x)$ and $f(x) \cdot g(x)$ in $\mathbb{Z}[x]$

(b) $f(x) + g(x)$ and $f(x) \cdot g(x)$ in $\mathbb{Z}_2[x]$

(c) $(f(x) \cdot g(x)) \cdot f(x)$ in $\mathbb{Z}[x]$

(d) $(f(x) \cdot g(x)) \cdot f(x)$ in $\mathbb{Z}_2[x]$

(e) $f(x) \cdot f(x) + f(x) \cdot g(x)$ in $\mathbb{Z}_2[x]$

2. Complete Example 16.3.3.

3. Prove that:

(a) The ring \mathbb{R} is a subring of the ring $\mathbb{R}[x]$.

(b) The ring $\mathbb{Z}[x]$ is a subring of the $\mathbb{Q}[x]$.

(c) The ring $\mathbb{Q}[x]$ is a subring of the ring $\mathbb{R}[x]$.

4. (a) Find all zeros of $x^4 + 1$ in \mathbb{Z}_3 . (b) Find all zeros of $x^5 + 1$ in \mathbb{Z}_5 .

5. Determine which of the following are reducible over \mathbb{Z}_2 . Explain.

(a) $f(x) = x^3 + 1$

(b) $g(x) = x^3 + x^2 + x$.

(c) $h(x) = x^3 + x^2 + 1$.

(d) $k(x) = x^4 + x^2 + 1$. (Be careful.)

6. Prove the second half of Theorem 16.3.3.

7. Give an example of the contention made in the last paragraph of this section.

8. Determine all zeros of $x^4 + 3x^3 + 2x + 4$ in $\mathbb{Z}_5[x]$

9. Show that $x^2 - 3$ is irreducible over \mathbb{Q} but reducible over the field of real numbers.

B Exercises

10. The definition of a vector space given in Chapter 13 holds over any field F , not just over the field of real numbers, where the elements of F are called scalars.

(a) Show that $F[x]$ is a vector space over F .

(b) Find a basis for $F[x]$ over F .

(c) What is the dimension of $F[x]$ over F ?

11. Prove Theorem 16.3.2.

(a) Show that the field \mathbb{R} of real numbers is a vector space over \mathbb{R} . Find a basis for this vector space. What is $\dim \mathbb{R}$ over \mathbb{R} ?

(b) Repeat part a for an arbitrary field F .

(c) Show that \mathbb{R} is a vector space over \mathbb{Q} .

16.4 Field Extensions

From high school algebra we realize that to solve a polynomial equation means to find its roots (or, equivalently, to find the zeros of the polynomials). From Example 16.3.5 of the previous section we know that the zeros may not lie in the given ground field. Hence, to solve a polynomial really involves two steps: first, find the zeros, and second, find the field in which the zeros lie. For economy's sake we would like this field to be the smallest field that contains all the zeros of the given polynomial. To illustrate this concept, let us reconsider Example 16.3.5.

Example 16.4.1. Let $f(x) = x^2 - 2 \in \mathbb{Q}[x]$. It is important to remember that we are considering $x^2 - 2$ over \mathbb{Q} , no other field. We would like to find all zeros of $f(x)$ and the smallest field, call it S for now, that contains them. The zeros are $x = \pm\sqrt{2}$, neither of which is an element of \mathbb{Q} . The set S we are looking for must satisfy the conditions:

- (1) S be a field.
- (2) S must contain \mathbb{Q} as a subfield,
- (3) S must contain all zeros of $f(x) = x^2 - 2$, and

By condition (3), $\sqrt{2}$ must be an element of S , and, if S is to be a field, the sum, product, difference, and quotient of elements in S must be in S . So $\sqrt{2}, (\sqrt{2})^2, (\sqrt{2})^3, \dots, \sqrt{2} + \sqrt{2}, \sqrt{2} - \sqrt{2}$, and $\sqrt{2}/\sqrt{2}$ must all be elements of S . Further, since S contains \mathbb{Q} as a subset, any element of \mathbb{Q} combined with $\sqrt{2}$ under any field operation must be an element of S . Hence, every element of the form $a + b\sqrt{2}$, where a and b can be any elements in \mathbb{Q} , is an element of S . We leave to the reader to show that S is a field (see Exercise 1 of this section). We note that the second zero of $x^2 - 2$, namely $-\sqrt{2}$, is an element of S . To see this, simply take $a = 0$ and $b = -1$. The field S is frequently denoted as $\mathbb{Q}(\sqrt{2})$, and it is referred to as an extension field of \mathbb{Q} . Note that the polynomial $x^2 - 2 = (x - \sqrt{2})(x + \sqrt{2})$ factors into linear factors, or *splits*, in $\mathbb{Q}(\sqrt{2})[x]$; that is, all coefficients of both factors are elements of the field $\mathbb{Q}(\sqrt{2})$.

Example 16.4.2. Consider the polynomial $g(x) = x^2 + x + 1 \in \mathbb{Z}_2[x]$. Let's repeat the previous example for $g(x)$ over \mathbb{Z}_2 . First, $g(0) = 1$ and $g(1) = 1$, so none of the elements of \mathbb{Z}_2 are zeros of $g(x)$. Hence, the zeros of $g(x)$ must lie in an extension field of \mathbb{Z}_2 . By Theorem 16.3.3, $g(x) = x^2 + x + 1$ can have at most two zeros. Let a be a zero of $g(x)$. Then the extension field S of \mathbb{Z}_2 must contain $a \cdot a = a^2, a^3, a + a, a + 1$, and so on. But, since $g(a) = 0$, we have $a^2 + a + 1 = 0$, or, equivalently, $a^2 = -(a + 1) = a + 1$ (remember, we are working in an extension of \mathbb{Z}_2). Note the recurrence relation.

So far our extension field S of \mathbb{Z}_2 is the set $\{0, 1, a, a + 1\}$. For S to be a field, all possible sums, products, differences, and quotients of elements in S must be in S . Let's try a few:

$$a + a = a(1 + 1) = a \cdot 0 = 0 \in S$$

Since $a + a = 0$, $-a = a$, which is in S . Adding three a 's together doesn't give us anything new: $a + a + a = a \in S$. In fact, na is in S for all possible positive integers n . Next,

$$\begin{aligned} a^3 &= a^2 \cdot a \\ &= (a + 1) \cdot a \\ &= a^2 + a \\ &= (a + 1) + a \\ &= 1 \in S \end{aligned}$$

Therefore, $a^{-1} = a + 1$ and $(a + 1)^{-1} = a$.

It is not difficult to see that a^n is in S for all positive n . Does S contain all zeros of $x^2 + x + 1$? Remember, $g(x)$ can have at most two distinct zeros and we called one of them a , so if there is a second, it must be $a + 1$. To see if $a + 1$ is indeed a zero of $g(x)$, simply compute $f(a + 1)$:

$$\begin{aligned} f(a + 1) &= (a + 1)^2 + (a + 1) + 1 \\ &= a^2 + 1 + a + 1 + 1 \\ &= a^2 + a + 1 \\ &= 0 \end{aligned}$$

Therefore, $a + 1$ is also a zero of $x^2 + x + 1$. Hence, $S = \{0, 1, a, a + 1\}$ is the smallest field that contains $\mathbb{Z}_2 = \{0, 1\}$ as a subfield and all zeros of $x^2 + x + 1$. This extension field is denoted by $\mathbb{Z}_2(a)$. Note that $x^2 + x + 1$ splits in $\mathbb{Z}_2(a)$; that is, it factors into linear factors in $\mathbb{Z}_2(a)$. We also observe that $\mathbb{Z}_2(a)$ is a field containing exactly four elements. By Theorem 16.2.4, we expected that $\mathbb{Z}_2(a)$ would be of order p^2 for some prime p and positive integer n . Also recall that all fields of order p^n are isomorphic. Hence, we have described all fields of order $2^2 = 4$ by finding the extension field of a polynomial that is irreducible over \mathbb{Z}_2 .

The reader might feel somewhat uncomfortable with the results obtained in Example 16.4.2. In particular, what is a ? Can we describe it through a known quantity? All we know about a is that it is a zero of $g(x)$ and that $a^2 = a + 1$. We could also say that $a(a + 1) = 1$, but we really expected more. However, should we expect more? In Example 16.4.1, $\sqrt{2}$ is a number we are more comfortable with, but all we really know

about it is that $\alpha = \sqrt{2}$ is the number such that $\alpha^2 = 2$. Similarly, the zero that the reader will obtain in Exercise 2 of this section is the imaginary number i . Here again, this is simply a symbol, and all we know about it is that $i^2 = -1$. Hence, the result obtained in Example 16.4.2 is not really that strange.

The reader should be aware that we have just scratched the surface in the development of topics in polynomial rings. One area of significant applications is in coding theory.

EXERCISES FOR SECTION 16.4

A Exercises

1. (a) Use the definition of a field to show that $\mathbb{Q}(\sqrt{2})$ is a field.
 (b) Use the definition of vector space to show that $\mathbb{Q}(\sqrt{2})$ is a vector space over \mathbb{Q} .
 (c) Prove that $\{1, \sqrt{2}\}$ is a basis for the vector space $\mathbb{Q}(\sqrt{2})$ over \mathbb{Q} , and, therefore, the dimension of $\mathbb{Q}(\sqrt{2})$ over \mathbb{Q} is 2.
2. (a) Determine the splitting field of $f(x) = x^2 + 1$ over \mathbb{R} . This means consider the polynomial $f(x) = x^2 + 1 \in \mathbb{R}[x]$ and find the smallest field that contains \mathbb{R} and all the zeros of $f(x)$. Denote this field by $\mathbb{R}(i)$.
 (b) $\mathbb{R}(i)$ is more commonly referred to by a different name. What is it?
 (c) Show that $\{1, i\}$ is a basis for the vector space $\mathbb{R}(i)$ over \mathbb{R} . What is the dimension of this vector space (over \mathbb{R})?
3. Determine the splitting field of $x^4 - 5x^2 + 6$ over \mathbb{Q} .
4. (a) Factor $x^2 + x + 1$ into linear factors in $\mathbb{Z}_2(a)$.
 (b) Write out the field tables for the field $\mathbb{Z}_2(a)$ and compare the results to the tables of Example 16.2.2.
 (c) Cite a theorem and use it to show why the results of part b were to be expected.
5. (a) Show that $x^3 + x + 1$ is irreducible over \mathbb{Z}_2 .
 (b) Determine the splitting field of $x^3 + x + 1$ over \mathbb{Z}_2 .
 (c) Use Theorem 16.2.4 to illustrate that you have described all fields of order 2^3 .
6. (a) List all polynomials of degree 1, 2, 3, and 4 over $\mathbb{Z}_2 = \text{GF}(2)$.
 (b) Use your results in part a and list all irreducible polynomials of degree 1, 2, 3, and 4.
 (c) Determine the splitting fields of each of the polynomials in part b.
 (d) What is the order of each of the splitting fields obtained in part c? Explain your results using Theorem 16.2.4.

16.5 Power Series

In Section 16.3 we found that a polynomial of degree n over a ring R is an expression of the form

$$f(x) = \sum_{i=0}^n a_i x^i = a_0 + a_1 x + a_2 x^2 + \cdots + a_n x^n, \quad n \geq 0,$$

where each of the a_i are elements of R and $a_n \neq 0$. In Section 8.5 we defined a generating function of a sequence s with terms s_0, s_1, s_2, \dots as the infinite sum

$$G(s, z) = \sum_{i=0}^{\infty} s_i z^i = s_0 + s_1 z + s_2 z^2 + \cdots$$

The main difference between these two expressions, disregarding notation, is that the latter is an infinite expression and the former is a finite expression. In this section we will extend the algebra of polynomials to the algebra of infinite expressions like $G(s, z)$, which are called *power series*.

Definition: Power Series. Let $[R; +, \cdot]$ be a ring. A power series over R is an expression of the form

$$f(x) = \sum_{i=0}^{\infty} a_i x^i = a_0 + a_1 x + a_2 x^2 + \cdots$$

where $a_1, a_2, a_3, \dots \in R$. The set of all such expressions is denoted by $R[[x]]$.

Our first observation in our comparison of $R[x]$ and $R[[x]]$ is that every polynomial is a power series and so $R[x] \subseteq R[[x]]$. This is true because a polynomial $a_0 + a_1 x + a_2 x^2 + \cdots + a_n x^n$ of degree n in $R[x]$, can be thought of as an infinite expression where $a_i = 0$ for $i > n$. In addition,

we will see that $R[[x]]$ is a ring with subring $R[x]$.

$R[[x]]$ is given a ring structure by defining addition and multiplication on power series as we did in $R[x]$, with the modification that, since we are dealing with infinite expressions, the sums and products will remain infinite expressions that we can determine term by term, as was done in Section 16.3.

Definition: Power Series Addition and Multiplication. Given power series

$$f(x) = \sum_{i=0}^{\infty} a_i x^i = a_0 + a_1 x + a_2 x^2 + \cdots$$

and

$$g(x) = \sum_{i=0}^{\infty} b_i x^i = b_0 + b_1 x + b_2 x^2 + \cdots$$

their sum is

$$f(x) + g(x) = \sum_{i=0}^{\infty} (a_i + b_i) x^i$$

and their product is

$$f(x) \cdot g(x) = \sum_{i=0}^{\infty} d_i x^i$$

where

$$d_i = \sum_{j=0}^i a_j b_{i-j}$$

Let's look at an example.

Example 16.5.1. (Example 8.5.3, Revisited.) Let

$$f(x) = \sum_{i=0}^{\infty} i x^i = 0 + 1x + 2x^2 + 3x^3 + \cdots$$

and

$$g(x) = \sum_{i=0}^{\infty} 2^i x^i = 1 + 2x + 4x^2 + 8x^3 + \cdots$$

be elements in $\mathbb{Z}[[x]]$. Let us compute $f(x) + g(x)$ and $f(x) \cdot g(x)$. First the sum:

$$\begin{aligned} f(x) + g(x) &= \sum_{i=0}^{\infty} i x^i + \sum_{i=0}^{\infty} 2^i x^i = \sum_{i=0}^{\infty} (i + 2^i) x^i \\ &= 1 + 3x + 6x^2 + 11x^3 + \cdots \end{aligned}$$

The product is a bit more involved:

$$\begin{aligned} f(x) \cdot g(x) &= \left(\sum_{i=0}^{\infty} i x^i \right) \left(\sum_{i=0}^{\infty} 2^i x^i \right) \\ &= (0 + 1x + 2x^2 + 3x^3 + \cdots)(1 + 2x + 4x^2 + 8x^3 + \cdots) \\ &= 0 \cdot 1 + (0 \cdot 2 + 1 \cdot 1)x + (0 \cdot 4 + 1 \cdot 2 + 2 \cdot 1)x^2 + \cdots \\ &= \sum_{i=0}^{\infty} d_i x^i \end{aligned}$$

where

$$d_i = \sum_{j=0}^i a_j b_{i-j} = \sum_{j=0}^i j 2^{i-j}$$

For example,

$$\begin{aligned} d_3 &= 0 \cdot 2^3 + 1 \cdot 2^2 + 2 \cdot 2^1 + 3 \cdot 2^0 \\ &= 0 + 4 + 4 + 3 \\ &= 11 \end{aligned}$$

Hence,

$$f(x) \cdot g(x) = x + 4x^2 + 11x^3 + \cdots$$

The First few terms of the product do not suggest a pattern but with *Mathematica*, we can get a closed form expression for the coefficients.

$$\text{Simplify}\left[\sum_{j=0}^i j 2^{i-j}\right]$$

$$-i + 2^{i+1} - 2$$

Therefore, $d_i = 2^{i+1} - i - 2$ and

$$\begin{aligned} f(x) \cdot g(x) &= \left(\sum_{i=0}^{\infty} i x^i\right) \left(\sum_{i=0}^{\infty} 2^i x^i\right) \\ &= \sum_{i=0}^{\infty} (2^{i+1} - i - 2) x^i \end{aligned}$$

We have shown that addition and multiplication in $R[[x]]$ is virtually identical to that in $R[x]$. The following theorem parallels Theorem 16.3.1, establishing the ring properties of $R[[x]]$.

Theorem 16.5.1. *Let $[R, +, \cdot]$ be a ring. Then:*

- (1) $R[[x]]$ is a ring under the operations of power series addition and multiplication, which depend on (are induced by) the operations in R .
- (2) If R is a commutative ring, then $R[[x]]$ is a commutative ring.
- (3) If R is a ring with unity, 1, then $R[[x]]$ is a ring with unity (the unity in $R[x]$ is $1 + 0x + 0x^2 + \dots$).
- (4) If R is an integral domain, then $R[[x]]$ is an integral domain.
- (5) If F is a field, then $F[[x]]$ is not a field. However, $F[[x]]$ is an integral domain.

We are most interested in the situation when the set of coefficients is a field. Theorem 16.5.1 indicates that when F is a field, $F[[x]]$ is an integral domain. A reason that $F[[x]]$ is not a field is the same as one that we can cite for $F[x]$, namely that x does not have multiplicative inverse in $F[[x]]$. With all of these similarities, one might wonder if the rings of polynomials and power series over a field are isomorphic. It turns out that they are not.

The difference between $F[x]$ and $F[[x]]$ become apparent when one studies which elements are units (i.e., elements that have multiplicative inverses) in each. First we prove that the only units in $F[x]$ are the nonzero constants — that is, the nonzero elements of F .

Theorem 16.5.2. *Let $[F; +, \cdot]$ be a field, $f(x)$ is a unit in $F[x]$ if and only if $f(x)$ is a nonzero element of F .*

Proof: (\Rightarrow) Let $f(x)$ be a unit in $F[x]$. Then $f(x)$ has a multiplicative inverse, call it $g(x)$, such that $f(x) \cdot g(x) = 1$. Hence, the $\deg(f(x) \cdot g(x)) = \deg(1) = 0$. But $\deg(f(x) \cdot g(x)) = \deg f(x) + \deg g(x)$. So $\deg f(x) + \deg g(x) = 0$, and since the degree of a polynomial is always nonnegative, this can only happen when the $\deg f(x) = \deg g(x) = 0$. Hence, $f(x)$ is a constant, an element of F , which is a unit if and only if it is nonzero.

(\Leftarrow) If $f(x)$ is a nonzero element of F , then it is a unit since F is a field. Thus it has an inverse, which is also in $F[x]$ and so $f(x)$ is a unit of $F[x]$. ■

Before we proceed to categorize the units in $F[[x]]$, we remind the reader that two power series $a_0 + a_1x + a_2x^2 + \dots$ and $b_0 + b_1x + b_2x^2 + \dots$ are equal if and only if corresponding coefficients are equal, $a_i = b_i$ for all $i \geq 0$.

Theorem 16.5.3. *Let $[F; +, \cdot]$ be a field. Then $f(x) = \sum_{i=0}^{\infty} a_i x^i$ is a unit of $F[[x]]$ if and only if $a_0 \neq 0$.*

Proof: (\Rightarrow) If $f(x)$ is a unit of $F[[x]]$, then there exists $g(x) = \sum_{i=0}^{\infty} b_i x^i$ in $F[[x]]$ such that

$$\begin{aligned} f(x) \cdot g(x) &= (a_0 + a_1x + a_2x^2 + \dots) \cdot (b_0 + b_1x + b_2x^2 + \dots) \\ &= 1 \\ &= 1 + 0x + 0x^2 + \dots \end{aligned}$$

Since corresponding coefficients in the equation above must be equal, $a_0 \cdot b_0 = 1$, which implies that $a_0 \neq 0$.

(\Leftarrow) Assume that $a_0 \neq 0$. To prove that $f(x)$ is a unit of $F[[x]]$ we need to find $g(x) = \sum_{i=0}^{\infty} b_i x^i$ in $F[[x]]$ such that

$$f(x) \cdot g(x) = \sum_{i=0}^{\infty} d_i x^i = 1.$$

If we use the formula for the coefficients d_0, d_1, d_2, \dots of $f(x) \cdot g(x)$ and equate coefficients, we will obtain

$$\begin{aligned}
d_0 &= a_0 \cdot b_0 = 1 \\
d_1 &= a_0 b_1 + a_1 b_0 = 0 \\
d_2 &= a_0 b_2 + a_1 b_1 + a_2 b_0 \\
&\vdots \\
d_s &= a_0 b_s + a_1 b_{s-1} + \cdots + a_s b_0 \\
&\vdots
\end{aligned}$$

Therefore, the existence of $g(x)$ is equivalent to the existence of a solution b_0, b_1, b_2, \dots , to the above system of equations. Since $a_0 \neq 0$, we can solve the first equation for b_0 . Then we can continue to the second equation and solve for b_1 , then b_2 can be found by solving the third equation, etc. Hence,

$$\begin{aligned}
b_0 &= a_0^{-1} \\
b_1 &= -a_0^{-1}(a_1 b_0) \\
b_2 &= -a_0^{-1}(a_1 b_1 + a_2 b_0) \\
&\vdots \\
b_s &= -a_0^{-1}(a_1 b_{s-1} + a_2 b_{s-2} + \cdots + a_s b_0) \\
&\vdots
\end{aligned}$$

Therefore the powers series $\sum_{i=0}^{\infty} b_i x^i$ is an expression whose coefficients lie in F and that satisfies the statement $f(x) \cdot g(x) = 1$. Hence, $g(x)$ is the multiplicative inverse of $f(x)$ and $f(x)$ is a unit..

Example 16.5.2. Let

$$\begin{aligned}
f(x) &= 1 + 2x + 3x^2 + 4x^3 + \cdots \\
&= \sum_{i=0}^{\infty} (i+1)x^i
\end{aligned}$$

be an element of $\mathbb{Q}[[x]]$. Then, by Theorem 16.5.3, since $a_0 = 1 \neq 0$, $f(x)$ is a unit and has an inverse, call it $g(x)$. To compute $g(x)$, we follow the procedure outlined in Theorem 16.5.3. Using the formulas for the b_i 's, we obtain

$$\begin{aligned}
b_0 &= 1 \\
b_1 &= -1(2 \cdot 1) = -2 \\
b_2 &= -1(2 \cdot (-2) + 3 \cdot 1) = 1 \\
b_3 &= -1(2 \cdot 1 + 3 \cdot (-2) + 4 \cdot 1) = 0 \\
b_4 &= -1(2 \cdot 0 + 3 \cdot 1 + 4 \cdot (-2) + 5 \cdot 1) = 0 \\
b_5 &= -1(2 \cdot 0 + 3 \cdot 0 + 4 \cdot (1) + 5 \cdot (-2) + 6 \cdot 1) = 0 \\
&\vdots \\
b_s &= -1(2 \cdot 0 + 3 \cdot 0 + \cdots + (s-2) \cdot 0 + (s-1) \cdot 1 + s \cdot (-2) + (s+1) \cdot 1) = 0 \quad s \geq 3
\end{aligned}$$

Hence, $g(x) = 1 - 2x + x^2$ is the multiplicative inverse of $f(x)$.

If we compare Theorems 16.5.2 and 16.5.3, we note that while the only elements in $F[x]$ that are units are the nonzero constants of F , the units in $F[[x]]$ are every single expression in x where $a_0 \neq 0$. So certainly $F[[x]]$ contains a wider variety of units than $F[x]$. Yet $F[[x]]$ is not a field, since $x \in F[[x]]$ is not a unit by Theorem 16.5.3. So concerning the algebraic structure of $F[[x]]$, we know that it is an integral domain that contains $F[x]$. If we allow our power series to take on negative exponents—that is, consider expressions of the form

$$f(x) = \sum_{i=-\infty}^{\infty} a_i x^i$$

where all but a finite number of terms with a negative index equal zero. These expressions are called *extended power series*. The set of all such expressions is a field, call it E . This set does contain, for example, the inverse of x namely x^{-1} . It can be shown that each nonzero element of E is a unit.

EXERCISES FOR SECTION 16.5

A Exercises

1. Let $f(x) = \sum_{i=0}^{\infty} a_i x^i$ and $g(x) = \sum_{i=0}^{\infty} b_i x^i$ be elements of $R[[x]]$. Let

$$f(x) \cdot g(x) = \sum_{i=0}^{\infty} d_i x^i = 1.$$

(a) Apply the distributive law repeatedly to

$$(a_0 + a_1 x + a_2 x^2 + \cdots) \cdot (b_0 + b_1 x + b_2 x^2 + \cdots)$$

to obtain the formula

$$d_s = \sum_{i=0}^s a_i b_{s-i}$$

for the coefficients of $f(x) \cdot g(x)$. Hence, you have shown that

$$f(x) \cdot g(x) = \sum_{s=0}^{\infty} \left(\sum_{i=0}^s a_i b_{s-i} \right) x^s$$

(b) Apply the above formula to the product in Example 16.5.1 and show that the result is the same as that obtained in this example.

2. (a) Prove that for any integral domain D , the following can be proven:

$$f(x) = \sum_{i=0}^{\infty} a_i x^i \text{ is a unit of } D[[x]] \text{ if and only if } a_0 \text{ is a unit in } D.$$

(b) Compare the statement in part a to that in Theorem 16.5.3.

(c) Give an example of the statement in part a in $\mathbb{Z}[[x]]$.

3. Use the formula for the product to verify that the expression $g(x)$ of Example 16.5.2 is indeed the inverse of $f(x)$.

4. (a) Determine the inverse of $f(x) = 1 + x + x^2 + \cdots = \sum_{i=0}^{\infty} x^i$ in $\mathbb{Q}[[x]]$.

(b) Repeat part a with $f(x)$ taken in $\mathbb{Z}_2[[x]]$.

(c) Use the method outlined in Chapter 8 to show that the power series $f(x) = \sum_{i=0}^{\infty} x^i$ is the rational generating function $\frac{1}{1-x}$. What is the inverse of this function? Compare your results with those in part a.

5. (a) Determine the inverse of $h(x) = \sum_{i=0}^{\infty} 2^i x^i$ in $\mathbb{Q}[[x]]$.

(b) Use the procedures in Chapter 8 to find a rational generating function for $h(x)$ in part a. Find the multiplicative inverse of this function.

6. Let $a(x) = 1 + 3x + 9x^2 + 27x^3 + \cdots = \sum_{i=0}^{\infty} 3^i x^i$ and

$$b(x) = 1 + x + x^2 + x^3 + \cdots = \sum_{i=0}^{\infty} x^i \text{ both in } \mathbb{R}[[x]].$$

(a) What are the first four terms (counting the constant term as the 0th term) of $a(x) + b(x)$?

(b) Find a closed form expression for $a(x)$.

(c) What are the first four terms of $a(x)b(x)$?

B Exercise

7. Write as an extended power series:

(a) $(x^4 - x^5)^{-1}$

(b) $(x^2 - 2x^3 + x^4)^{-1}$

SUPPLEMENTARY EXERCISES FOR CHAPTER 16

Section 16.1

- Expand $(A + B)^2$ in the ring $[M_{n \times n}(\mathbb{R}); +, \cdot]$.
 - Will your result be similar for any noncommutative ring?
- Expand $(A + B)^3$ in the ring $[M_{n \times n}(\mathbb{R}); +, \cdot]$.
 - Will your result be similar for any noncommutative ring?
- Let D be the set of all 2×2 diagonal matrices over the real numbers.
 - Prove that D is a subring of $[M_{2 \times 2}(\mathbb{R}); +, \cdot]$, hence a ring under the usual operations.
 - Prove that D is a commutative ring with unity.
 - Is the cancellation law true in D ?
- Use the definition of a ring to convince yourself that $R = \{a + b\sqrt{2} \mid a, b \in \mathbb{Z}\}$ is a ring. A common name given this ring is $\mathbb{Z}[\sqrt{2}]$.
 - What is the unity of $\mathbb{Z}[\sqrt{2}]$?
 - Prove that $\mathbb{Z}[\sqrt{2}]$ is an integral domain.
- It can be shown, in general, that if R is any ring, $[M_{n \times n}(R); +, \cdot]$ is a ring.
 - How many elements are there in the ring $R = [M_{2 \times 2}(\mathbb{Z}_2); +, \cdot]$? What are the zero and unity of R ?
 - Determine all solutions of the equation $X^2 - I = 0$ in R .
- Find all six units of $[M_{2 \times 2}(\mathbb{Z}_2); +, \cdot]$. Hint: The set of units is closed with respect to multiplication and one of them is $\begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}$.
- Let $A = \left\{ \begin{pmatrix} a & 0 \\ 0 & 0 \end{pmatrix} \mid a \in \mathbb{R} \right\}$ then A is a ring under matrix addition and multiplication. Prove that A is isomorphic to the ring of real numbers.

Section 16.2

- Show that \mathbb{Z}_2 is a subfield of the field given in Example 16.2.2, or equivalently, that the field in this example is an extension field of \mathbb{Z}_2 .
- Show that a and b are the two roots of the equation $x^2 + x + 1 = 0$ in the field of Example 16.2.2.
- Let $A = \left\{ \begin{pmatrix} a & b \\ -b & a \end{pmatrix} \mid a, b \in \mathbb{R} \right\}$. Prove that A with matrix addition and multiplication is isomorphic to the ring of complex numbers, \mathbb{C} .

Section 16.3

- Find all rational zeros (roots) of $f(x) = x^4 - 6x^3 + 10x^2 - 6x + 9$ and factor $f(x)$ into irreducible factors in $\mathbb{Q}[x]$.
- Determine all zeros of $f(x) = x^3 + 1$ in the field of Example 16.2.2, and express $f(x)$ as a product of irreducible factors over that field.
- Repeat Exercise 12 for $g(x) = x^2 + x^2 + x$,
- Find all five roots of $f(x) = x^3 + 7x$ in \mathbb{Z}_8 . Explain why this does not contradict Theorem 16.3.4.

Exercises 15 to 20 develop an introduction to polynomial codes. In Chapter 15 we introduced group codes. Here, we will discuss another code that uses polynomials. A k -tuple in \mathbb{Z}_2^k can be identified with a polynomial of degree $k - 1$ in the integral domain $\mathbb{Z}_2[x]$ and conversely. We do this by associating a k -tuple with the coefficients of a polynomial starting with the constant term. For example, the 5-tuple $(1, 0, 1, 1, 0)$ is viewed as the polynomial $1 + 0x + 1x^2 + 1x^3 + 0x^4 = 1 + x^2 + x^3$. If we define addition and multiplication on \mathbb{Z}_2^k based on polynomial operations, we will have highly structured codes. For the actual code, we present an example where $k = 7$.

15. To add k -tuples, we can take two equivalent approaches. We can either simply add the k -tuples coordinatewise as in any direct product, or we can convert the k -tuples to polynomials of degree $k - 1$ or less, add them, and then write down the coefficients of the sum.

(a) For each of the following pairs of add and multiply the pairs k -tuples, where k varies, compute their sum. Use both ways to add for at least one part.

(i) $(0, 1, 0)$ and $(1, 1, 1)$

(ii) $(0, 1, 0, 1)$ and $(1, 1, 0, 1)$

(iii) $(1, 1, 1, 0, 1, 0, 1)$ and $(0, 0, 0, 0, 1, 0, 0)$

(iv) $(1, 0, 0, 1, 1, 1, 1)$ and $(0, 0, 0, 1, 0, 0, 0)$

(b) What relationship between polynomials and k -tuples makes it possible to do this addition two different ways to get the same sum.

16. The encoding of a string of bits is based on polynomial division. Given a four bit message, we make the bits coefficients of a sixth degree polynomial, $b_3x^3 + b_4x^4 + b_5x^5 + b_6x^6$ which we can also express in \mathbb{Z}_2^6 as $(0, 0, 0, b_3, b_4, b_5, b_6)$, we divide this polynomial by $p(x) = 1 + x + x^3$ and add the remainder to the “message polynomial”. The quotient in the division is discarded. Thus, if the remainder, which must be a polynomial of degree less than 2, is $b_0 + b_1x + b_2x^2$, the encoded message is the string of bits $(b_0, b_1, b_2, b_3, b_4, b_5, b_6)$.

(a) Encode the following elements of \mathbb{Z}_2^6 as described above.

(a) $(0, 0, 0, 1, 1, 0, 1)$

(b) $(0, 0, 0, 1, 1, 1, 1)$

(c) $(0, 0, 0, 0, 0, 1, 0)$

(b) Prove that the encoded message will always represent a polynomial which is evenly divisible by the polynomial $p(x)$ that is used to encode the message.

17. A single bit error in the transmission of our seven bit encoded message $(b_0, b_1, b_2, b_3, b_4, b_5, b_6)$ can be thought of as a monomial expression x^j , where $0 \leq j \leq 6$. The effect of an error occurring is to add that monomial to the encoded message. So if the last bit is transmitted incorrectly, the monomial x^6 is added and the received bit sequence is $(b_0, b_1, b_2, b_3, b_4, b_5, b_6 + 1)$. If no error takes place, we can think of the zero polynomial being added. Prove that if an error takes place, the received bit string represents a polynomial that is *not* a multiple of $p(x)$.

18. There are seven different single bit errors. Let's focus on what happens if an error occurs in the last bit. If the error occurs in the last bit and the received bit string represents the polynomial $m(x)$, show that the remainder upon dividing $m(x)$ by $p(x)$ will be the same for all possible values of $m(x)$. What is that remainder? This is called the *syndrome* for an error in the last bit.

19. What are the syndromes for each of the other error positions? Let's agree to number them 0th through 6th, so the 6th position syndrome was determined above. What is the syndrome if no error occurs?

20. Assuming no more than a single bit error in the transmission of seven bits, what is the transmitted bit string, given these received strings?

(a) $(0, 1, 0, 1, 0, 1, 1)$

(b) $(1, 1, 1, 0, 0, 0, 0)$

(c) $(0, 0, 1, 1, 0, 1, 0)$

Section 16.4

21. In Exercise 5 of Section 16.4 you constructed GF(8) using $x^3 + x + 1$. Show that GF(8) can also be obtained by using the polynomial $g(x) = x^3 + x^2 + 1$.

22. (a) Show that $f(x) = x^4 + x + 1$ is irreducible over \mathbb{Z}_2 .

(b) Describe the splitting field of $f(x)$ over \mathbb{Z}_2 .

(c) Let a be a zero of $f(x)$. Show that each nonzero element of the splitting field in part (b) can be described as a power of a .

Section 16.5

23. Review Example 16.5.2. Derive the multiplicative inverse of $1 - 2x + x^2$ by doing repeated polynomial division, as suggested by the following first step:

$$\begin{array}{r}
 1 \\
 1 - 2x + x^2 \overline{) 1} \\
 \underline{1 - 2x + x^2} \\
 2x - x^2
 \end{array}$$

24. Use polynomial long division to obtain the power series representation of $\frac{1}{1+x^3}$ over \mathbb{Q} . What is the inverse of the power series you obtained?

25. Find the generating function for the sequence defined by the difference equation $a_k = a_{k-1} + a_{k-2}$, $k \geq 2$, with $a_0 = a_1 = 1$ the indicated fields.

(a) \mathbb{Q}

(b) \mathbb{Z}_2

(c) \mathbb{Z}_3

26. Determine the inverse of each of the power series in Exercise 25.

27. Recall from high school algebra that any quadratic with real coefficients, of the form $ax^2 + bx + c = 0$, $a \neq 0$, can be solved using the quadratic formula:

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

(a) Does this formula always produce zeros in \mathbb{R} ?

(b) Use the quadratic formula to solve $x^2 + x + 1 = 0$ in \mathbb{Z}_3 ,

(c) Use the quadratic formula to solve $x^2 + 2 = 0$ in \mathbb{Z}_3 .

(d) Use the quadratic formula to solve $x^2 + x + 2 = 0$ in \mathbb{Z}_3 .

(e) What observations do you have based on parts (b) - (d)?

Solutions to Odd-numbered Exercises

CHAPTER 11

Section 11.1

1. (a) Commutative, and associative. Notice that zero is the identity for addition, but it is not a positive integer.)

(b) Commutative, associative, and has an identity (1)

(c) Commutative, associative, has an identity (1), and is idempotent

(d) Commutative, associative, and idempotent

(e) None. Note: $2 @ (3 @ 3) = 512$
 $(2 @ 3) @ 3 = 64$

and while $a @ 1 = a$, $1 @ a = 1$.

3. $a, b \in A \cap B \Rightarrow a, b \in A$ by the definition of intersection
 $\Rightarrow a * b \in A$ by the closure of A with respect to $*$

Similarly, $a, b \in A \cap B \Rightarrow a * b \in B$. Therefore, $a * b \in A \cap B$.

The set of positive integers is closed under addition, and so is the set of negative integers, but $1 + -1 = 0$. Therefore, their union, the nonzero integers, is not closed under addition.

5. Let \mathbb{N} be the set of all nonnegative integers (the natural numbers).

(a) $*$ is commutative since $|a - b| = |b - a|$ for all $a, b \in \mathbb{N}$

(b) $*$ is not associative. Take $a = 1, b = 2$, and $c = 3$, then

$$(a * b) * c = ||1 - 2| - 3| = 2, \text{ and}$$

$$a * (b * c) = |1 - |2 - 3|| = 0.$$

(c) Zero is the identity for $*$ on \mathbb{N} , since

$$a * 0 = |a + 0| = a = |0 - a| = 0 * a.$$

(d) $a^{-1} = a$ for each $a \in \mathbb{N}$, since

$$a * a = |a - a| = 0.$$

(e) $*$ is not idempotent, since, for $a \neq 0$,

$$a * a = 0 \neq a.$$

Section 11.2

1. The terms "generic" and "trade" for prescription drugs are analogous to "generic" and "concrete" algebraic systems. Generic aspirin, for example, has no name, whereas Bayer, Tylenol, Bufferin, and Anacin are all trade or specific types of aspirins. The same can be said of a generic group $[G, *]$ where G is a nonempty set and $*$ is a binary operation on G . When examples of typical domain elements can be given along with descriptions of how operations act on them, such as \mathbb{Q}^* or $M_{2 \times 2}(\mathbb{R})$, then the system is concrete (has a specific name, as with the aspirin). Generic is a way to describe a general algebraic system, whereas a concrete system has a name or symbols making it distinguishable from other systems.

3. b, d, e, and f.

5. (a) $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$, abelian

	I	R_1	R_2	F_1	F_2	F_3
I	I	R_1	R_2	F_1	F_2	F_3
R_1	R_1	R_2	I	F_2	F_3	F_1
R_2	R_2	I	R_1	F_3	F_1	F_2
F_1	F_1	F_2	F_3	I	R_2	R_1
F_2	F_2	F_1	F_3	R_1	I	R_2
F_3	F_3	F_2	F_1	R_2	R_1	I

This group is non-abelian since, for example, $F_1 F_2 = R_2$ and $F_2 F_1 = R_2$.

(c) $4! = 24, n!$

7. The identity is e . $a * b = c, a * c = b, b * c = a$, and $[V, *]$ is abelian. (This group is commonly called the Klein-4 group.)

Section 11.3

1. (a) f is injective: $f(x) = f(y) \Rightarrow a * x = a * y$
 $\Rightarrow x = y$ (by left cancellation)

f is surjective: For all b , $f(x) = b$ has the solution $a^{-1} * b$.

(b) Functions of the form $f(x) = a + x$, where a is any integer, are bijections

3. Basis: ($n = 2$) $(a_1 * a_2)^{-1} = a_2^{-1} * a_1^{-1}$ by Theorem 11.3.4.

Induction: Assume that for some $n \geq 2$,

$$(a_1 * a_2 * \cdots * a_n)^{-1} = a_n^{-1} * \cdots * a_2^{-1} * a_1^{-1}$$

We must show that

$$(a_1 * a_2 * \cdots * a_n * a_{n+1})^{-1} = a_{n+1}^{-1} * a_n^{-1} * \cdots * a_2^{-1} * a_1^{-1}$$

This can be accomplished as follows:

$$\begin{aligned} (a_1 * a_2 * \cdots * a_n * a_{n+1})^{-1} &= ((a_1 * a_2 * \cdots * a_n) * a_{n+1})^{-1} \text{ by the associative law} \\ &= a_{n+1}^{-1} * (a_1 * a_2 * \cdots * a_n)^{-1} \text{ by the basis} \\ &= a_{n+1}^{-1} * (a_n^{-1} * \cdots * a_2^{-1} * a_1^{-1}) \text{ by the induction hypothesis} \\ &= a_{n+1}^{-1} * a_n^{-1} * \cdots * a_2^{-1} * a_1^{-1} \text{ by the associative law} \quad \blacksquare \end{aligned}$$

5. (a) Let $p(n)$ be, where a is any element of group $[G; *]$. First we will prove that $p(n)$ is true for all $n \geq 0$.

First, we would need to prove a lemma that we leave to the reader, that if $n \geq 0$, and a is any group element, $a * a^n = a^n * a$.

Basis: If $n = 0$, Using the definition of the zero exponent, $(a^0)^{-1} = e^{-1} = e$, while $(a^{-1})^0 = e$. Therefore, $p(0)$ is true.

Induction: Assume that for some $n \geq 0$, $p(n)$ is true.

$$\begin{aligned} (a^{n+1})^{-1} &= (a^n * a)^{-1} \text{ by the definition of exponentiation} \\ &= a^{-1} * (a^n)^{-1} \text{ by Theorem 11.3.4} \\ &= a^{-1} * (a^{-1})^n \text{ by the induction hypothesis} \\ &= (a^{-1})^{n+1} \text{ by the lemma} \end{aligned}$$

If n is negative, then $-n$ is positive and

$$\begin{aligned} a^{-n} &= (((a^{-1})^{-1})^{-n}) \\ &= (a^{-1})^{-(n)} \text{ since the property is true for positive numbers} \\ &= (a^{-1})^n \end{aligned}$$

(b) For $m > 1$, let $p(m)$ be $a^{n+m} = a^n * a^m$ for all $n \geq 1$. The basis for this proof follows directly from the basis for the definition of exponentiation.

Induction: Assume that for some $m > 1$, $p(m)$ is true. Then

$$\begin{aligned} a^{n+(m+1)} &= a^{(n+m)+1} \text{ by the associativity of integer addition} \\ &= a^{n+m} * a^1 \text{ by the definition of exponentiation} \\ &= (a^n * a^m) * a^1 \text{ by the induction hypothesis} \\ &= a^n * (a^m * a^1) \text{ by associativity} \\ &= a^n * a^{m+1} \text{ by the definition of exponentiation} \end{aligned}$$

(c) Let $p(m)$ be $(a^n)^m = a^{n*m}$ for all integers n .

Basis: $(a^m)^0 = e$ and $a^{m*0} = a^0 = e$ therefore, $p(0)$ is true.

Induction: Assume that $p(m)$ is true for some $m > 0$,

$$\begin{aligned} (a^n)^{m+1} &= (a^n)^m * a^n \text{ definition of exponentiation} \\ &= a^{n*m} * a^n \text{ by the induction hypothesis} \\ &= a^{n*m+n} \text{ by part (a) of this problem} \\ &= a^{n(m+1)} \end{aligned}$$

Finally, if m is negative, we can verify that $(a^n)^m = a^{nm}$ using many of the same steps as the proof of part (a).

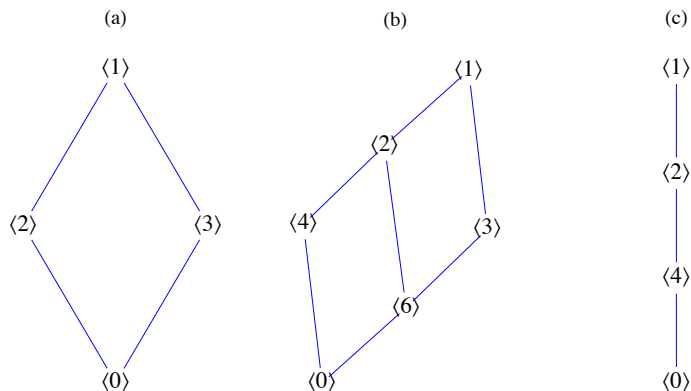
Section 11.4

1. (a) 2 (b) 5 (c) 0
 (d) 0 (e) 2 (f) 2
 (g) 1 (h) 3
3. (a) 1 (b) 1 (c) $m(4) = r(4)$, where $m = 11q + r, 0 \leq r < 11$.
5. Since the solutions, if they exist, must come from \mathbb{Z}_2 , substitution is the easiest approach.
 - (a) 1 is the only solution, since $1^2 +_2 1 = 0$ and $0^2 +_2 1 = 1$
 - (b) No solutions, since $0^2 +_2 0 +_2 1 = 1$, and $1^2 +_2 1 +_2 1 = 1$

7. Hint: Prove by induction on m that you can divide any positive integer into m . That is, let $p(m)$ be "For all n greater than zero, there exist unique integers q and r such that $n = mq + r$." In the induction step, divide n into $m - n$.

Section 11.5

1. a and c
3. $\{I, R_1, R_2\}, \{I, F_1\}, \{I, F_2\}$, and $\{I, F_3\}$ are all the proper subgroups of R_3 .
5. (a) $\langle 1 \rangle = \langle 5 \rangle = \mathbb{Z}_6$
 - $\langle 2 \rangle = \langle 4 \rangle = \{2, 4, 0\}$
 - $\langle 3 \rangle = \{3, 0\}$
 - $\langle 0 \rangle = \{0\}$
- (b) $\langle 1 \rangle = \langle 5 \rangle = \langle 7 \rangle = \langle 11 \rangle = \mathbb{Z}_{12}$
 - $\langle 2 \rangle = \langle 10 \rangle = \{2, 4, 6, 8, 10, 0\}$
 - $\langle 3 \rangle = \langle 9 \rangle = \{3, 6, 9, 0\}$
 - $\langle 4 \rangle = \langle 8 \rangle = \{4, 8, 0\}$
 - $\langle 6 \rangle = \{6, 0\}$
 - $\langle 0 \rangle = \{0\}$
- (c) $\langle 1 \rangle = \langle 3 \rangle = \langle 5 \rangle = \langle 7 \rangle = \mathbb{Z}_8$
 - $\langle 2 \rangle = \langle 6 \rangle = \{2, 4, 6, 0\}$
 - $\langle 4 \rangle = \{4, 0\}$
 - $\langle 0 \rangle = \{0\}$



(d) Based on the ordering diagrams in parts a through c, we would expect to see an ordering diagram similar to the one for divisors on $\{1, 2, 3, 4, 6, 8, 12, 24\}$ (the divisors of 24) if we were to examine the subgroups of \mathbb{Z}_{24} . This is indeed the case.

7. Assume that H and K are subgroups of group G , and that, as in Figure 11.5.1, there are elements $x \in H - K$ and $y \in K - H$. Consider the product $x * y$. Where could it be placed in the Venn diagram? If we can prove that it must lie in the outer region, $H^c \cap K^c = (H \cup K)^c$, then we have proven that $H \cup K$ is not closed under $*$ and can't be a subgroup of G . Assume that $x * y \in H$. Since x is in H , x^{-1} is in H and so by closure

$$x^{-1} * (x * y) = y \in H$$

which is a contradiction. Similarly, $x * y \notin K$. ■

One way to interpret this theorem is that no group is the union of two groups.

Section 11.6

1. Table of $\mathbb{Z}_2 \times \mathbb{Z}_3$:

		y					
x	*	{0, 0}	{0, 1}	{0, 2}	{1, 0}	{1, 1}	{1, 2}
	{0, 0}	{0, 0}	{0, 1}	{0, 2}	{1, 0}	{1, 1}	{1, 2}
	{0, 1}	{0, 1}	{0, 2}	{0, 0}	{1, 1}	{1, 2}	{1, 0}
	{0, 2}	{0, 2}	{0, 0}	{0, 1}	{1, 2}	{1, 0}	{1, 1}
	{1, 0}	{1, 0}	{1, 1}	{1, 2}	{0, 0}	{0, 1}	{0, 2}
	{1, 1}	{1, 1}	{1, 2}	{1, 0}	{0, 1}	{0, 2}	{0, 0}
	{1, 2}	{1, 2}	{1, 0}	{1, 1}	{0, 2}	{0, 0}	{0, 1}

The only two proper subgroups are $\{(0, 0), (1, 0)\}$ and $\{(0, 0), (0, 1), (0, 2)\}$

3. (a) (i) $a + b$ could be $(1, 0)$ or $(0, 1)$.

(ii) $a + b = (1, 1)$.

(b) (i) $a + b$ could be $(1, 0, 0)$, $(0, 1, 0)$, or $(0, 0, 1)$.

(ii) $a + b = (1, 1, 1)$.

(c) (i) $a + b$ has exactly one 1.

(ii) $a + b$ has all 1's.

5. (a) No, 0 is not an element of $\mathbb{Z} \times \mathbb{Z}$.

(b) Yes.

(c) No, $(0, 0)$ is not an element of this set.

(d) No, the set is not closed: $(1, 1) + (2, 4) = (3, 5)$ and $(3, 5)$ is not in the set.

(e) Yes.

Section 11.7

1. (a) Yes, $f(n, x) = (x, n)$ for $(n, x) \in \mathbb{Z} \times \mathbb{R}$ is an isomorphism.

(b) No, $\mathbb{Z}_2 \times \mathbb{Z}$ has a finite proper subgroup while $\mathbb{Z} \times \mathbb{Z}$ does not.

(c) No.

(d) Yes.

(e) No.

(f) Yes, one isomorphism is defined by $f(a_1, a_2, a_3, a_4) = \begin{pmatrix} a_1 & a_2 \\ a_3 & a_4 \end{pmatrix}$.

(g) Yes, one isomorphism is defined by $f(a_1, a_2) = (a_1, 10^{a_2})$.

(h) Yes.

(i) Yes $f(k) = k(1, 1)$.

3. Consider 3 groups G_1 , G_2 , and G_3 with operations $*$, \diamond , and \square , respectively.. We want to show that if G_1 is isomorphic to G_2 , and if G_2 is isomorphic to G_3 , then G_1 is isomorphic to G_3 .

G_1 isomorphic to $G_2 \Rightarrow$ there exists an isomorphism $f : G_1 \rightarrow G_2$

G_2 isomorphic to $G_3 \Rightarrow$ there exists an isomorphism $g : G_2 \rightarrow G_3$

If we compose g with f , we get the function $g \circ f : G_1 \rightarrow G_3$. By Theorems 7.3.2 and 7.3.3, $g \circ f$ is a bijection, and if $a, b \in G_1$,

$$\begin{aligned}(g \circ f)(a * b) &= g(f(a * b)) \\ &= g(f(a) \diamond f(b)) \text{ since } f \text{ is an isomorphism} \\ &= g(f(a)) \square g(f(b)) \text{ since } g \text{ is an isomorphism} \\ &= (g \circ f)(a) * (g \circ f)(b)\end{aligned}$$

Therefore, $g \circ f$ is an isomorphism from G_1 into G_3 , proving that "is isomorphic to" is transitive.

5. \mathbb{Z}_8 , $\mathbb{Z}_2 \times \mathbb{Z}_4$, and \mathbb{Z}_2^3 . One other is the fourth dihedral group, introduced in Section 15.3.

7. Let G be an infinite cyclic group generated by a . Then, using multiplicative notation, $G = \{a^n \mid n \in \mathbb{Z}\}$.

The map $T : G \rightarrow \mathbb{Z}$ defined by $T(a^n) = n$ is an isomorphism. This is indeed a function, since $a^n = a^m$ implies $n = m$. Otherwise, a would have a finite order and would not generate G .

(a) T is one-to-one, since $T(a^n) = T(a^m)$ implies $n = m$, so $a^n = a^m$.

(b) T is onto, since for any $n \in \mathbb{Z}$, $T(a^n) = n$.

$$\begin{aligned}\text{(c)} \quad T(a^n * a^m) &= T(a^{n+m}) \\ &= n + m \\ &= T(a^n) + T(a^m)\end{aligned}$$

Supplementary Exercises—Chapter 11

1. (a) With respect to V under $+$, the identity is a ; and $-a = a$, $-b = c$, and $-c = b$.

(b) With respect to V under \cdot , the identity is b . Inverses: $b^{-1} = b$, $c^{-1} = c$, and a has no inverse,

(c) \cdot is distributive over $+$ since $x \cdot (y + z) = x \cdot y + x \cdot z$ for each of the 27 ways that the variables x, y , and z can be assigned values from V . However, $+$ is not distributive over \cdot since $b + (a \cdot c) = b$, while $(b + a) \cdot (b + c) = a$,

3. By Theorem 7.3.4 every bijection has an inverse, so \circ has the inverse property on S . If $f \in S$,

$$f \circ f^{-1} = f^{-1} \circ f = i \Rightarrow f \text{ inverts } f^{-1}, \text{ or } (f^{-1})^{-1} = f.$$

Therefore, inversion of functions has the involution property.

5. If a and b are odd integers, $a = 2j + 1$ and $b = 2k + 1$ for $j, k \in \mathbb{Z}$. $ab = (2j + 1)(2k + 1) = 2(2jk + j + k) + 1$, which is an odd integer. Since 1 is odd and $1 + 1$ is even, the odds are not closed under addition. The even integers are closed under both addition and multiplication. If a and b are even, $a = 2j$ and $b = 2k$ for some $j, k \in \mathbb{Z}$, $a + b = 2j + 2k = 2(j + k)$, which is even, and $ab = (2j)(2k) = 2(2jk)$, which is also even.

7. That $\text{GL}(2, \mathbb{R})$ is a group follows from laws of matrix algebra. In addition to being associative, matrix multiplication on two-by-two matrices has an identity I , and if $A \in \text{GL}(2, \mathbb{R})$, it has an inverse by the definition of $\text{GL}(2, \mathbb{R})$. The inverse of A is in $\text{GL}(2, \mathbb{R})$ since it has an inverse: $(A^{-1})^{-1} = A$.

9. If $a, b, c \in \mathbb{R}$,

$$\begin{aligned}(a * b) * c &= (a + b + 5) * c \\ &= a + b + 5 + c + 5 \\ &= a + b + c + 10\end{aligned}$$

$a * (b * c)$ is also equal to $a + b + c + 10$, and so $*$ is associative. To find the identity we solve $a * e = a$ for e :

$$a * e = a \Rightarrow a + e + 5 = a \Rightarrow e = -5.$$

If a is a real number, the inverse of a is determined by solving the equation $a * x = -5$;

$$a * x = -5 \Rightarrow a + x + 5 = -5 \Rightarrow x = -a - 10$$

Since a is real, $-a - 10$ is real, and so $*$ has the inverse property.

11. By Supplementary Exercise 2 of this chapter, the identity for $*$ is 2 and $*$ is associative. All that is left to show is that $*$ has the inverse

property. If $a \in \mathbb{Q}^+$, $a * x = 2 \Rightarrow x = \frac{4}{a}$; hence $a^{-1} = \frac{4}{a}$, which is also a positive rational number.

13. Recall that matrix multiplication is the operation on $\text{GL}(2, \mathbb{R})$.

$$\begin{aligned} A X B = C &\Rightarrow X B = A^{-1} C \quad (\text{multiply on the left by } A^{-1}) \\ &\Rightarrow X = A^{-1} C B^{-1} \quad (\text{multiply on the right by } B^{-1}) \end{aligned}$$

$$X = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{3} \end{pmatrix} \begin{pmatrix} 2 & 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{2} & -\frac{1}{2} \\ -\frac{1}{2} & 1 \end{pmatrix} = \begin{pmatrix} \frac{1}{4} & 0 \\ -\frac{1}{6} & \frac{1}{3} \end{pmatrix}$$

15. (a) 1 (b) 4 (c) 0 (d) 3

17. (a) $\langle 1 \rangle = \{1\}$, $\langle 3 \rangle = \{1, 3\}$, $\langle 5 \rangle = \{1, 5\}$, and $\langle 7 \rangle = \{1, 7\}$.

(b) No, because no cyclic subgroup equals $U(\mathbb{Z}_8)$.

19. (a) $A, B \in \text{SL}(2, \mathbb{R}) \Rightarrow |A| = |B| = 1$.

$$\begin{aligned} |AB| &= |A| \cdot |B| = 1 \cdot 1 = 1 \Rightarrow AB \in \text{SL}(2, \mathbb{R}) \\ &\Rightarrow \text{SL}(2, \mathbb{R}) \text{ is closed with respect to matrix multiplication} \end{aligned}$$

(b) $|I| = 1 \Rightarrow I \in \text{SL}(2, \mathbb{R})$

(c) $A \in \text{SL}(2, \mathbb{R}) \Rightarrow |A| = 1$

$$|A^{-1}| = |A|^{-1} = 1 \Rightarrow A^{-1} \in \text{SL}(2, \mathbb{R})$$

21. Yes, S is a submonoid of $B_{3 \times 3}$. The zero matrix is in S since it is the matrix of the empty relation, which is symmetric. Furthermore, if A and B are matrices of symmetric relations,

$$\begin{aligned} (A + B)_{ij} &= A_{ij} + B_{ij} \quad \text{definition of matrix addition} \\ &= A_{ji} + B_{ji} \quad \text{since both } A \text{ and } B \text{ are symmetric} \\ &= (A + B)_{ji} \quad \text{definition of matrix addition} \end{aligned}$$

Therefore, $A + B$ is symmetric, which means that it is the matrix of a symmetric relation and that relation is in S .

23. (a) $(1, 4, 20)$ (b) $(-1, 0, -1)$ (c) $(1/3, 4)$ (d) $(-2, -3, -5)$

25. The groups in parts a and c are abelian, since each factor is abelian. The group in part b is non-abelian, since one of its factors, $\text{GL}(2, \mathbb{R})$, is non-abelian.

27. Since $\langle 4 \rangle = \{0, 4, 8, 12\}$ is a cyclic group and has order four, it must be isomorphic to \mathbb{Z}_4 .

29. (a) There exists a "dictionary" that allows us to translate between the two systems in such a way that any true fact in one is translated to a true fact in the other.

(b) If one system is familiar to you, the other one should be familiar too.

(c) If $(p \wedge \neg q) \Leftrightarrow 0$, and $(p \wedge q) \Leftrightarrow 0$, then $p \Leftrightarrow 0$.

31. The key to this exercise is to identify the fact that adding two complex numbers entails adding two pairs of numbers, the real and imaginary parts. If we simply rename these parts the first and second parts, then we are doing \mathbb{R}^2 addition. This suggests the function $T: \mathbb{C} \rightarrow \mathbb{R}^2$ where $T(a + bi) = (a, b)$. For any two complex numbers $a + bi$ and $c + di$,

$$\begin{aligned} T((a + bi) + (c + di)) &= T((a + c) + (b + d)i) \quad \text{definition of } + \text{ in } \mathbb{C} \\ &= (a + c, b + d) \quad \text{definition of } T \\ &= (a, b) + (c, d) \quad \text{definition of } + \text{ in } \mathbb{R}^2 \\ &= T(a + bi) + T(c + di) \quad \text{definition of } T \end{aligned}$$

Since T has an inverse ($T^{-1}(a, b) = a + bi$), T is an isomorphism and so the two groups are isomorphic.

It should be noted that T is not the only isomorphism between these two groups. For example $U(a + bi) = (b, a)$ defines an isomorphism.

33. The key here is to realize that both groups consist of elements that are constructed from four real numbers and that you operate on elements by adding four different pairs of real numbers. An isomorphism from \mathbb{R}^4 into $M_{2 \times 2}(\mathbb{R})$ is

$$T(a, b, c, d) = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

There are an infinite number of isomorphism in this case. This one is the most obvious.

CHAPTER 12

Section 12.1

1. (a) $\{(4/3, 1/3)\}$

(b) $\{(-3 - 0.5x_3, 11 - 4x_3, x_3) \mid x_3\}$

(c) $\{(-5, 14/5, 8/5)\}$

(d) $\{(6.25 - 2.5x_3, -0.75 + 0.5x_3, x_3) \mid x_3 \in \mathbb{R}\}$

3. (a) $\{(1.2, 2.6, 4.5)\}$

(b) $\{(-6x_3 + 5, 2x_3 + 1, x_3) \mid x_3 \in \mathbb{R}\}$

(c) $\{(-9x_3 + 3, 4, x_3) \mid x_3 \in \mathbb{R}\}$

(d) $\{(3x_4 + 1, -2x_4 + 2, x_4 + 1, x_4) \mid x_4 \in \mathbb{R}\}$

5. (a) $\{(3, 0)\}$

$$\begin{aligned}
 \text{(b)} \quad & \begin{pmatrix} 1 & 1 & 2 & 1 \\ 1 & 2 & 4 & 4 \\ 1 & 3 & 3 & 0 \end{pmatrix} \xrightarrow{\substack{-R_1 + R_2 \\ -R_1 + R_3}} \begin{pmatrix} 1 & 1 & 2 & 1 \\ 0 & 1 & 2 & 3 \\ 0 & 2 & 1 & -1 \end{pmatrix} \\
 & \xrightarrow{\substack{-R_2 + R_1 \\ -2R_2 + R_3}} \begin{pmatrix} 1 & 0 & 0 & -2 \\ 0 & 1 & 2 & 3 \\ 0 & 0 & -3 & -7 \end{pmatrix} \\
 & \xrightarrow{\substack{-\frac{1}{3}R_3}} \begin{pmatrix} 1 & 0 & 0 & -2 \\ 0 & 1 & 2 & 3 \\ 0 & 0 & 1 & \frac{7}{3} \end{pmatrix} \\
 & \xrightarrow{\substack{-\frac{1}{2}R_3 + R_2}} \begin{pmatrix} 1 & 0 & 0 & -2 \\ 0 & 1 & 2 & 3 \\ 0 & 0 & 1 & \frac{7}{3} \end{pmatrix}
 \end{aligned}$$

The row reduction can be done with *Mathematica*:

$$\begin{aligned}
 & \text{RowReduce}\left[\begin{pmatrix} 1 & 1 & 2 & 1 \\ 1 & 2 & 4 & 4 \\ 1 & 3 & 3 & 0 \end{pmatrix}\right] \\
 & \begin{pmatrix} 1 & 0 & 0 & -2 \\ 0 & 1 & 0 & -\frac{5}{3} \\ 0 & 0 & 1 & \frac{7}{3} \end{pmatrix}
 \end{aligned}$$

In any case, the solution set is $\{(-2, -5/3, 7/3)\}$

7. Proof: Since b is the $n \times 1$ matrix of 0's, let's call it $\mathbf{0}$. Let S be the set of solutions to $AX = \mathbf{0}$. If X_1 and X_2 be in S . Then

$$A(X_1 + X_2) = AX_1 + AX_2 = \mathbf{0} + \mathbf{0} = \mathbf{0}$$

so $X_1 + X_2 \in S$; that is, S is closed under addition.

The identity of \mathbb{R}^n is $\mathbf{0}$, which is in S . Finally, let X be in S . Then

$$A(-X) = -(AX) = -\mathbf{0} = \mathbf{0},$$

and so $-X$ is also in S .

Section 12.2

(a) $\begin{pmatrix} \frac{15}{11} & \frac{30}{11} \\ \frac{3}{11} & -\frac{5}{11} \end{pmatrix}$

$$(b) \begin{pmatrix} -20 & \frac{21}{2} & \frac{9}{2} & -\frac{3}{2} \\ 2 & -1 & 0 & 0 \\ -4 & 2 & 1 & 0 \\ 7 & -\frac{7}{2} & -\frac{3}{2} & \frac{1}{2} \end{pmatrix}$$

(c) The inverse does not exist. When the augmented matrix is row-reduced (see below), the last row of the first half cannot be manipulated to match the identity matrix.

$$(d) \begin{pmatrix} 1 & 0 & 0 \\ -3 & 1 & 1 \\ -4 & 1 & 2 \end{pmatrix}$$

(e) The inverse does not exist.

$$(f) \begin{pmatrix} 9 & -36 & 30 \\ -36 & 192 & -180 \\ 30 & -180 & 180 \end{pmatrix}$$

5. The solutions are in the solution section of Section 12.1, exercise 1. We illustrate with the outline of the solution to Exercise 1c of Section 12.1.

$$\begin{pmatrix} 1 & 1 & 2 \\ 1 & 2 & -1 \\ 1 & 3 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ 5 \end{pmatrix}$$

$$A^{-1} = \begin{pmatrix} 1 & 1 & 2 \\ 1 & 2 & -1 \\ 1 & 3 & 1 \end{pmatrix}^{-1} = \frac{1}{5} \begin{pmatrix} 5 & 5 & -5 \\ -2 & -1 & 3 \\ 1 & -2 & 1 \end{pmatrix}$$

$$\text{and } \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = A^{-1} \begin{pmatrix} 1 \\ -1 \\ 5 \end{pmatrix} = \begin{pmatrix} -5 \\ \frac{14}{5} \\ \frac{8}{5} \end{pmatrix}$$

Section 12.3

3. (b) Yes

7. If the matrices are named B , A_1 , A_2 , A_3 , and A_4 , then

$$B = \frac{8}{3} A_1 + \frac{5}{3} A_2 + \frac{-5}{3} A_3 + \frac{23}{3} A_4.$$

9. (a) If $x_1 = (1, 0)$, $x_2 = (0, 1)$, and $y = (b_1, b_2)$, then

$$y = b_1 x_1 + b_2 x_2.$$

If $x_1 = (3, 2)$, $x_2 = (2, 1)$, and $y = (b_1, b_2)$, then

$$y = (-b_1 + 2b_2)x_1 + (2b_1 - 3b_2)x_2.$$

The second linear combination can be computed using *Mathematica* as follows.

Solve[$\mathbf{c}_1 \{3, 2\} + \mathbf{c}_2 \{2, 1\} == \{\mathbf{b}_1, \mathbf{b}_2\}, \{\mathbf{c}_1, \mathbf{c}_2\}$]

$\{\{c_1 \rightarrow 2b_2 - b_1, c_2 \rightarrow 2b_1 - 3b_2\}\}$

(b) If $y = (b_1, b_2)$ is any vector in \mathbb{R}^2 , then

$$y = (-3b_1 + 4b_2)x_1 + (-b_1 + b_2)x_2 + (0)x_3$$

(c) One solution is to add any vector(s) to x_1 , x_2 , and x_3 of part b.

(d) $2, n$

(e) If the matrices are A_1 , A_2 , A_3 , and A_4 , then

$$\begin{pmatrix} x & y \\ z & w \end{pmatrix} = x A_1 z + y A_2 + z A_3 + w A_4$$

(f) $a_0 + a_1 x + a_2 x^2 + a_3 x^3 = a_0(1) + a_1(x) + a_2(x^2) + a_3(x^3).$

11. (a) The set is linearly independent: let a and b be scalars such that $a(4, 1) + b(1, 3) = (0, 0)$, then

$$\begin{aligned} 4a + b &= 0 \quad \text{and} \\ a + 3b &= 0 \end{aligned}$$

which has $a = b = 0$ as its only solutions. The set generates all of \mathbb{R}^2 : let (a, b) be an arbitrary vector in \mathbb{R}^2 . We want to show that we can always find scalars β_1 and β_2 such that $\beta_1(4, 1) + \beta_2(1, 3) = (a, b)$. This is equivalent to finding scalars such that $4\beta_1 + \beta_2 = a$ and $\beta_1 + 3\beta_2 = b$. This system has a unique solution $\beta_1 = \frac{3a-b}{11}$, and $\beta_2 = \frac{4b-a}{11}$. Therefore, the set generates \mathbb{R}^2 .

13. (d) They are isomorphic. Once you have completed part (a) of this exercise, the following translation rules will give you the answer to parts (b) and (c),

$$(a, b, c, d) \leftrightarrow \begin{pmatrix} a & b \\ c & d \end{pmatrix} \leftrightarrow a + bx + cx^2 + dx^2$$

Section 12.4

1. (a) Any nonzero multiple of $\begin{pmatrix} 1 \\ -1 \end{pmatrix}$ is an eigenvector associated with $\lambda = 1$.

(b) Any nonzero multiple of $\begin{pmatrix} 1 \\ 2 \end{pmatrix}$ is an eigenvector associated with $\lambda = 4$.

(c) Let $x_1 = \begin{pmatrix} a \\ -a \end{pmatrix}$ and $x_2 = \begin{pmatrix} b \\ 2b \end{pmatrix}$. You can verify that $c_1 x_1 + c_2 x_2 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ if and only if $c_1 = c_2 = 0$. Therefore, $\{x_1, x_2\}$ is linearly independent.

3. (c) You should obtain $\begin{pmatrix} 4 & 0 \\ 0 & 1 \end{pmatrix}$ or $\begin{pmatrix} 1 & 0 \\ 0 & 4 \end{pmatrix}$, depending on how you order the eigenvalues.

5. (a) If $P = \begin{pmatrix} 2 & 1 \\ 3 & -1 \end{pmatrix}$, then $P^{-1}AP = \begin{pmatrix} 4 & 0 \\ 0 & -1 \end{pmatrix}$.

(b) If $P = \begin{pmatrix} 1 & 1 \\ 7 & 1 \end{pmatrix}$, then $P^{-1}AP = \begin{pmatrix} 5 & 0 \\ 0 & -1 \end{pmatrix}$.

(c) If $P = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$, then $P^{-1}AP = \begin{pmatrix} 3 & 0 \\ 0 & 4 \end{pmatrix}$.

(d) If $P = \begin{pmatrix} 1 & -1 & 1 \\ -1 & 4 & 2 \\ -1 & 1 & 1 \end{pmatrix}$, then $P^{-1}AP = \begin{pmatrix} -2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$.

(e) A is not diagonalizable. Five is a double root of the characteristic equation, but has an eigenspace with dimension only 1.

(f) If $P = \begin{pmatrix} 1 & 1 & 1 \\ -2 & 0 & 1 \\ 1 & -1 & 1 \end{pmatrix}$, then $P^{-1}AP = \begin{pmatrix} 3 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$.

7. (b) This is a direct application of the definition of matrix multiplication. Let $A_{(i)}$ stand for the i^{th} row of A , and let $P^{(j)}$ stand for the j^{th} column of P . Hence the j^{th} column of the product AP is

$$\begin{pmatrix} A_{(1)}P^{(j)} \\ A_{(2)}P^{(j)} \\ \vdots \\ A_{(n)}P^{(j)} \end{pmatrix}$$

Hence, $(AP)^{(j)} = A(P^{(j)})$ for $j = 1, 2, \dots, n$. Thus, each column of AP depends on A and the j^{th} column of P .

Section 12.5

3. If we introduce the superfluous equation $1 = 0 \cdot S_{k-1} + 1$ we have the system

$$\begin{aligned} S_k &= 5S_{k-1} + 4 \\ 1 &= 0 \cdot S_{k-1} + 1 \end{aligned}$$

which, in matrix form, is:

$$\begin{aligned}
 \begin{pmatrix} S_k \\ 1 \end{pmatrix} &= \begin{pmatrix} 5 & 4 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} S_{k-1} \\ 1 \end{pmatrix} \\
 &= \begin{pmatrix} 5 & 4 \\ 0 & 1 \end{pmatrix}^k \begin{pmatrix} S_0 \\ 1 \end{pmatrix} \\
 &= \begin{pmatrix} 5 & 4 \\ 0 & 1 \end{pmatrix}^k \begin{pmatrix} 0 \\ 1 \end{pmatrix}
 \end{aligned}$$

Let $A = \begin{pmatrix} 5 & 4 \\ 0 & 1 \end{pmatrix}$. We want to diagonalize A ; that is, find a matrix P such that $P^{-1} A P = D$, where D is a diagonal matrix, or

$$A = P D P^{-1} \Rightarrow A^k = P D^k P^{-1}$$

Diagonalizing A :

$$|A - cI| = \begin{vmatrix} 5-c & 4 \\ 0 & 1-c \end{vmatrix} = (5-c)(1-c)$$

The eigenvalues are $c = 1$ and $c = 5$. If $c = 1$,

$$\begin{pmatrix} 4 & 4 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

which implies $x_1 + x_2 = 0$, or $x_2 = -x_1$, and so $\begin{pmatrix} 1 \\ -1 \end{pmatrix}$ is an eigenvector associated with 1.

If $c = 5$,

$$\begin{pmatrix} 0 & 4 \\ 0 & -4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Rightarrow x_2 = 0.$$

Therefore, $\begin{pmatrix} 1 \\ 0 \end{pmatrix}$ is an eigenvector associated with 5. Combining the two eigenvectors, we get

$$A = \begin{pmatrix} 1 & 1 \\ -1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 5 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ -1 & 0 \end{pmatrix}^{-1} = \begin{pmatrix} 1 & 1 \\ -1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 5 \end{pmatrix} \begin{pmatrix} 0 & -1 \\ 1 & 1 \end{pmatrix}$$

and

$$\begin{aligned}
 A^k &= \begin{pmatrix} 1 & 1 \\ -1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 5 \end{pmatrix}^k \begin{pmatrix} 0 & -1 \\ 1 & 1 \end{pmatrix} \\
 &= \begin{pmatrix} 1 & 1 \\ -1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 5^k \end{pmatrix} \begin{pmatrix} 0 & -1 \\ 1 & 1 \end{pmatrix} \\
 &= \begin{pmatrix} 5^k & 5^k - 1 \\ 0 & 1 \end{pmatrix}
 \end{aligned}$$

Hence, $\begin{pmatrix} S_k \\ 1 \end{pmatrix} = \begin{pmatrix} 5^k & 5^k - 1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 5^k - 1 \\ 1 \end{pmatrix}$ and finally, $S_k = 5^k - 1$.

5. Since $A = A^1 = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix}$, there are 0 paths of length 1 from: node c to node a, node b to node h, and node a to node c; and there is 1 path of length 1 for every other pair of nodes.

(b) The characteristic polynomial is

$$|A - cI| = \begin{vmatrix} 1-c & 1 & 0 \\ 1 & -c & 1 \\ 0 & 1 & 1-c \end{vmatrix} = -c^3 + 2c^2 + c - 2$$

Solving the characteristic equation $-c^3 + 2c^2 + c - 2 = 0$ we find solutions 1, 2, and -1.

If $c = 1$, we find the associated eigenvector by finding a nonzero solution to

$$\begin{pmatrix} 0 & 1 & 0 \\ 1 & -1 & 1 \\ 0 & 1 & 0 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

One of these, which will be the first column of P , is $\begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}$

If $c = 2$, the system $\begin{pmatrix} -1 & 1 & 0 \\ 1 & -2 & 1 \\ 0 & 1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$ yields eigenvectors, including $\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$, which will be the second column of P .

If $c = -1$, then the system determining the eigenvectors is

$$\begin{pmatrix} 2 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

and we can select $\begin{pmatrix} 1 \\ -2 \\ 1 \end{pmatrix}$, although any nonzero multiple of this vector could be the third column of P .

(c) Assembling the results of (b) we have $P = \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & -2 \\ -1 & 1 & 1 \end{pmatrix}$.

$$\begin{aligned} A^4 &= P \begin{pmatrix} 1^4 & 0 & 0 \\ 0 & 2^4 & 0 \\ 0 & 0 & (-1)^4 \end{pmatrix} P^{-1} = P \begin{pmatrix} 1 & 0 & 0 \\ 0 & 16 & 0 \\ 0 & 0 & 1 \end{pmatrix} P^{-1} \\ &= \begin{pmatrix} 1 & 16 & 1 \\ 0 & 16 & -2 \\ -1 & 16 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{2} & 0 & -\frac{1}{2} \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ \frac{1}{6} & -\frac{1}{3} & \frac{1}{6} \end{pmatrix} \\ &= \begin{pmatrix} 6 & 5 & 5 \\ 5 & 6 & 5 \\ 5 & 5 & 6 \end{pmatrix} \end{aligned}$$

Hence there are five different paths of length 4 between distinct vertices, and six different paths that start and end at the same vertex. The reader can verify these facts from Figure 12.4.1.

$$7. (a) e^A = \begin{pmatrix} e & e \\ 0 & 0 \end{pmatrix}, e^B = \begin{pmatrix} 0 & 0 \\ 0 & e^2 \end{pmatrix}, \text{ and } e^{A+B} = \begin{pmatrix} e & e^2 - e \\ 0 & e^2 \end{pmatrix}$$

(b) Let $\mathbf{0}$ be the zero matrix, $e^{\mathbf{0}} = I + \mathbf{0} + \frac{\mathbf{0}^2}{2} + \frac{\mathbf{0}^3}{6} + \dots = I$.

(c) Assume that A and B commute. We will examine the first few terms in the product $e^A e^B$. The pattern that is established does continue in general. In what follows, it is important that $AB = BA$. For example, in the last step, $(A+B)^2$ expands to $A^2 + AB + BA + B^2$, not $A^2 + 2AB + B^2$, if we can't assume commutativity.

$$\begin{aligned} e^A e^B &= \left(\sum_{k=0}^{\infty} \frac{A^k}{k!} \right) \left(\sum_{k=0}^{\infty} \frac{B^k}{k!} \right) \\ &= \left(I + A + \frac{A^2}{2} + \frac{A^3}{6} + \dots \right) \left(I + B + \frac{B^2}{2} + \frac{B^3}{6} + \dots \right) \\ &= I + A + B + \frac{A^2}{2} + AB + \frac{B^2}{2} + \frac{A^3}{6} + \frac{A^2 B}{2} + \frac{A B^2}{2} + \frac{B^3}{6} + \dots \\ &= I + (A+B) + \frac{1}{2} (A^2 + 2AB + B^2) + \frac{1}{6} (A^3 + 3A^2 B + 3A B^2 + B^3) + \dots \\ &= I + (A+B) + \frac{1}{2} (A+B)^2 + \frac{1}{6} (A+B)^3 + \dots \\ &= e^{A+B} \end{aligned}$$

$$e^A e^B = \left(\sum_{k=0}^{\infty} \frac{A^k}{k!} \right) \cdot \left(\sum_{k=0}^{\infty} \frac{B^k}{k!} \right)$$

$(A+B)^2 \text{ for } 2AB + A^2 + B^2$

(d) Since A and $-A$ commute, we can apply part d;

$$\begin{aligned}
 e^A e^{-A} &= e^{A+(-A)} \\
 &= e^{\mathbf{0}} \\
 &= I \quad \text{by part } b \text{ of this problem.}
 \end{aligned}$$

Supplementary Exercises—Chapter 12

1. (a) $x_1 = x_2 = x_3 = 1$

(b) $x_1 = \frac{1}{2}, x_2 = 0, x_3 = \frac{1}{2}$

3.
$$\begin{pmatrix} -8 & -4 & 1 \\ 7 & 3 & -1 \\ -5 & -2 & 1 \end{pmatrix}$$

5. Suppose that A^{-1} exists and that $\alpha_1(Ax_1) + \alpha_2(Ax_2)$ is equal to the zero vector, $\mathbf{0}$. By applying several laws of matrix algebra, this implies that

$$\begin{aligned}
 A(\alpha_1 x_1 + \alpha_2 x_2) &= \mathbf{0} \Rightarrow \alpha_1 x_1 + \alpha_2 x_2 = \mathbf{0} \quad \text{since } A^{-1} \text{ exists} \\
 &\Rightarrow \alpha_1 = \alpha_2 = 0 \quad \text{since } \{x_1, x_2\} \text{ is a basis} \\
 &\Rightarrow \{Ax_1, Ax_2\} \text{ is linearly independent}
 \end{aligned}$$

To see that $\{Ax_1, Ax_2\}$ also spans \mathbb{R}^2 , let $b \in \mathbb{R}^2$, we note that since $\{x_1, x_2\}$ is a basis, it will span $A^{-1}b$:

$$\alpha_1 x_1 + \alpha_2 x_2 = A^{-1}b \quad \text{for some } \alpha_1, \alpha_2 \in \mathbb{R}.$$

Using laws of matrix algebra:

$$\begin{aligned}
 \alpha_1 (Ax_1) + \alpha_2 (Ax_2) &= A(\alpha_1 x_1 + \alpha_2 x_2) \\
 &= A(A^{-1}b) \\
 &= b
 \end{aligned}$$

Hence, b is a linear combination of Ax_1 and Ax_2 .

If A has no inverse, then $Ax = \mathbf{0}$ has a nonzero solution y , which is spanned by the vectors x_1 and x_2 : $y = \alpha_1 x_1 + \alpha_2 x_2$, where not both of the α 's are zero.

$$\begin{aligned}
 Ay = \mathbf{0} &\Rightarrow A(\alpha_1 x_1 + \alpha_2 x_2) = \mathbf{0} \\
 &\Rightarrow \alpha_1 (Ax_1) + \alpha_2 (Ax_2) = \mathbf{0} \\
 &\Rightarrow \{Ax_1, Ax_2\} \text{ is linearly dependent}
 \end{aligned}$$

7. (b) $-X = X$

(c) $2^6 = 64$, since each entry can take on two possible values.

9. $A = P^{-1}DP \Rightarrow A^{100} = P^{-1}D^{100}P$

$$\begin{pmatrix} 0.6 & 0.2 \\ 0.4 & 0.8 \end{pmatrix} = \frac{1}{3} \begin{pmatrix} 1 & 2 \\ 2 & -1 \end{pmatrix} \begin{pmatrix} 1^{100} & 0 \\ 0 & 0.4^{100} \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 2 & -1 \end{pmatrix} \approx \begin{pmatrix} \frac{1}{3} & \frac{1}{3} \\ \frac{2}{3} & \frac{2}{3} \end{pmatrix}$$

Note: $0.4^{100} = 1.60694 \times 10^{-40} \approx 0$.

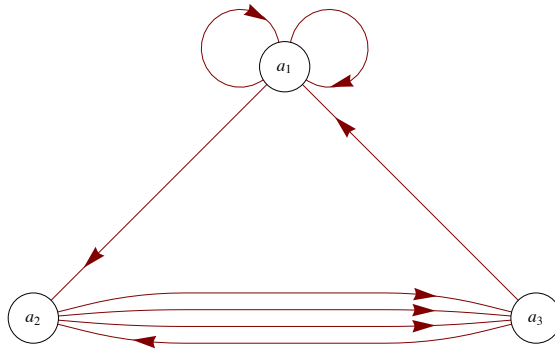
11. (a) $\lambda = 0, \pm\sqrt{2}$

(b)
$$B = PDP^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 0 & -2 \end{pmatrix} \begin{pmatrix} 4 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & 2 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & \frac{1}{2} \\ 0 & 0 & -\frac{1}{2} \end{pmatrix}$$

13. (a) Let the vertices be a_1, a_2 , and a_3 ; and use the convenient matrix representation

$$\begin{matrix} & a_1 & a_2 & a_3 \\ \begin{matrix} a_1 \\ a_2 \\ a_3 \end{matrix} & \begin{pmatrix} 2 & 1 & 0 \\ 1 & 0 & 3 \\ 1 & 1 & 0 \end{pmatrix} \end{matrix}$$

one sees immediately, for example, that there are 3 different edges from a_2 to a_3 , so that the multigraph is



(b) $A^2 = \begin{pmatrix} 5 & 2 & 3 \\ 5 & 4 & 0 \\ 3 & 1 & 3 \end{pmatrix}$ and by Theorem 12.5.1, $(A^2)_{ij}$ is the number of paths of length 2 from a_i to a_j . For example, the reader can verify from the graph that there are 3 different paths of length 2 from a_1 to a_3 .

CHAPTER 13

Section 13.1

1. (a) 1, 5 (b) 5

(c) 30 (d) 30

(e) See Figure 13.4.1 with $0 = 1, a_1 = 2, a_2 = 3, a_3 = 5, b_1 = 6, b_2 = 10, b_3 = 15$, and $1 = 30$

3. Solution for Hasse diagram (b):

(a)

lub	a_1	a_2	a_3	a_4	a_5
a_1	a_1	a_2	a_3	a_4	a_5
a_2	a_2	a_2	a_4	a_4	a_5
a_3	a_3	a_4	a_3	a_4	a_5
a_4	a_4	a_4	a_4	a_4	a_5
a_5	a_5	a_5	a_5	a_5	a_5

glb	a_1	a_2	a_3	a_4	a_5
a_1	a_1	a_1	a_1	a_1	a_1
a_2	a_1	a_2	a_1	a_2	a_2
a_3	a_1	a_1	a_3	a_3	a_3
a_4	a_1	a_2	a_3	a_4	a_4
a_5	a_1	a_2	a_3	a_4	a_5

(b) a_1 is the least element and a_5 is the greatest element.

Partial solution for Hasse diagram (f):

(a) $\text{lub}(a_2, a_3)$ and $\text{lub}(a_4, a_5)$ do not exist.

(b) No greatest element exists, but a_1 is the least element.

5. If 0 and $0'$ are distinct least elements, then

$$\left. \begin{array}{l} 0 \leq 0' \text{ since } 0 \text{ is a least element} \\ 0' \leq 0 \text{ since } 0' \text{ is a least element} \end{array} \right\} \Rightarrow 0 = 0' \text{ by antisymmetry, a contradiction. } \blacksquare$$

Section 13.2

1. Assume to the contrary that a and b have two different greatest lower bounds, and call them g and h . Then $g \geq h$ since g is a greatest lower bound and $h \geq g$ since h is a greatest lower bound. Therefore, by antisymmetry $h = g$.

3. (a) See Table 13.3.1 for the statements of these laws. Most of the proofs follow from the definition of gcd and lcm.

(b) (partial) We prove two laws as examples.

Commutative law of join: Let $[L, \vee, \wedge]$ be a lattice, $a, b \in L$. We must prove that $a \vee b = b \vee a$.

Proof: By the definition of least upper bound, $a \vee b \geq b$ and $a \vee b \geq a$ therefore, by Exercise 4, part c, $a \vee b \geq b \vee a$. Similarly, $b \vee a \geq a \vee b$, and by antisymmetry $a \vee b = b \vee a$. \blacksquare

Idempotent law (for join): We must prove that for all $a \in L$, $a \vee a = a$.

Proof: By the reflexive property of \leq , $a \leq a$ and hence, by 4(c), $a \leq a \vee a$. But a is an upper bound for a ; hence $a \geq a \vee a$. By antisymmetry, $a = a \vee a$. ■

Section 13.3

1.

B	Complement of B
\emptyset	A
$\{a\}$	$\{b, c\}$
$\{b\}$	$\{a, c\}$
$\{c\}$	$\{a, b\}$
$\{a, b\}$	$\{c\}$
$\{a, c\}$	$\{b\}$
$\{b, c\}$	$\{a\}$
A	\emptyset

This lattice is a Boolean algebra since it is a distributive complemented lattice.

3. a and g.

5. (a) $S^* : a \vee b = a$ if $a \geq b$

(b) $S : A \cap B = A$ if $A \subseteq B$

$S^* : A \cup B = A$ if $A \supseteq B$

(c) Yes

(d) $S : p \wedge q \Leftrightarrow p$ if $p \Rightarrow q$

$S^* : p \vee q \Leftrightarrow p$ if $q \Rightarrow p$

(e) Yes

7. **Definition: Boolean Algebra Isomorphism.** $[B, \wedge, \vee, -]$ is isomorphic to $[B', \wedge, \vee, -]$ if and only if there exists a function $T : B \rightarrow B'$ such that

(a) T is a bijection;

(b) $T(a \wedge b) = T(a) \wedge T(b)$ for all $a, b \in B$

(c) $T(a \vee b) = T(a) \vee T(b)$ for all $a, b \in B$

(d) $T(\bar{a}) = \widetilde{T(a)}$ for all $a \in B$.

Section 13.4

1. (a) For $a = 3$ we must show that for each $x \in D_{30}$ one of the following is true: $x \wedge 3 = 3$ or $x \wedge 3 = 1$. We do this through the following table:

x	verification
1	$1 \wedge 3 = 1$
2	$2 \wedge 3 = 1$
3	$3 \wedge 3 = 3$
5	$5 \wedge 3 = 1$
6	$6 \wedge 3 = 3$
10	$10 \wedge 3 = 1$
15	$15 \wedge 3 = 3$
30	$30 \wedge 3 = 3$

For $a = 5$, a similar verification can be performed.

(b) $6 = 2 \vee 3$, $10 = 2 \vee 5$, $15 = 3 \vee 5$, and $30 = 2 \vee 3 \vee 5$.

3. If $B = D_{30}$ then $A = \{2, 3, 5\}$ and D_{30} is isomorphic to $\mathcal{P}(A)$, where

$1 \leftrightarrow \emptyset$	$5 \leftrightarrow \{5\}$	
$2 \leftrightarrow \{2\}$	$10 \leftrightarrow \{2, 5\}$	and
$3 \leftrightarrow \{3\}$	$15 \leftrightarrow \{3, 5\}$	Join \leftrightarrow Union
$6 \leftrightarrow \{2, 3\}$	$30 \leftrightarrow \{2, 3, 5\}$	Meet \leftrightarrow Intersection
		Complement \leftrightarrow Set Complement

5. Assume that $x \neq 0$ or 1 is the third element of a Boolean algebra. Then there is only one possible set of tables for join and meet, all following from required properties of the Boolean algebra.

\vee	0	x	1
0	0	x	1
x	x	x	1
1	1	1	1

\wedge	0	x	1
0	0	0	0
x	0	x	x
1	0	x	1

Next, to find the complement of x we want y such that $x \wedge y = 0$ and $x \vee y = 1$. No element satisfies both conditions; hence the lattice is not complemented and cannot be a Boolean algebra. The lack of a complement can also be seen from the ordering diagram from which \wedge and \vee must be derived.

7. Let X be any countably infinite set, such as the integers. A subset of X is *cofinite* if it is finite or its complement is finite. The set of all cofinite subsets of X is:

(a) Countably infinite - this might not be obvious, but here is a hint. Assume $X = \{x_0, x_1, x_2, \dots\}$. For each finite subset A of X , map that set to the integer

$$\sum_{i=0}^{\infty} \chi_A(x_i) 2^i$$

You can do a similar thing to sets that have a finite complement, but map them to negative integers. Only one minor adjustment needs to be made to accommodate both the empty set and X .

(b) Closed under union

(c) Closed under intersection, and

(d) Closed under complementation.

Therefore, if $B = \{A \subseteq X : A \text{ is cofinite}\}$, then B is a countable Boolean algebra under the usual set operations.

Section 13.5

1. (a)

\vee	(0, 0)	(0, 1)	(1, 0)	(1, 1)
(0, 0)	(0, 0)	(0, 1)	(1, 0)	(1, 1)
(0, 1)	(0, 1)	(0, 1)	(1, 1)	(1, 1)
(1, 0)	(1, 0)	(1, 1)	(1, 0)	(1, 1)
(1, 1)	(1, 1)	(1, 1)	(1, 1)	(1, 1)

\wedge	(0, 0)	(0, 1)	(1, 0)	(1, 1)
(0, 0)	(0, 0)	(0, 0)	(0, 0)	(0, 0)
(0, 1)	(0, 0)	(0, 1)	(0, 0)	(0, 1)
(1, 0)	(0, 0)	(0, 0)	(1, 0)	(1, 0)
(1, 1)	(0, 0)	(0, 1)	(1, 0)	(1, 1)

u	\bar{u}
(0, 0)	(1, 1)
(0, 1)	(1, 0)
(1, 0)	(0, 1)
(1, 1)	(0, 0)

(b) The graphs are isomorphic.

(c) (0, 1) and (1, 0)

3. (a) (1, 0, 0, 0), (0, 1, 0, 0), (0, 0, 1, 0), and (0, 0, 0, 1) are the atoms.

(b) The n -tuples of 0's and 1's with exactly one 1.

Section 13.6

1 (a)

$$M_1(x_1, x_2) = 0$$

$$M_2(x_1, x_2) = (\bar{x}_1 \wedge \bar{x}_2)$$

$$M_3(x_1, x_2) = (\bar{x}_1 \wedge x_2)$$

$$M_4(x_1, x_2) = (x_1 \wedge \bar{x}_2)$$

$$M_5(x_1, x_2) = (x_1 \wedge x_2)$$

$$M_6(x_1, x_2) = ((\bar{x}_1 \wedge \bar{x}_2) \vee (\bar{x}_1 \wedge x_2)) = \bar{x}_1$$

$$M_7(x_1, x_2) = ((\bar{x}_1 \wedge \bar{x}_2) \vee (x_1 \wedge \bar{x}_2)) = \bar{x}_2$$

$$M_8(x_1, x_2) = ((\bar{x}_1 \wedge \bar{x}_2) \vee (x_1 \wedge x_2)) = ((x_1 \wedge x_2) \vee (\bar{x}_1 \wedge \bar{x}_2))$$

$$M_9(x_1, x_2) = ((\bar{x}_1 \wedge x_2) \vee (x_1 \wedge \bar{x}_2)) = ((x_1 \wedge \bar{x}_2) \vee (\bar{x}_1 \wedge x_2))$$

$$M_{10}(x_1, x_2) = ((\bar{x}_1 \wedge x_2) \vee (x_1 \wedge x_2)) = x_2$$

$$M_{11}(x_1, x_2) = ((x_1 \wedge \bar{x}_2) \vee (x_1 \wedge x_2)) = x_1$$

$$M_{12}(x_1, x_2) = ((\bar{x}_1 \wedge \bar{x}_2) \vee (\bar{x}_1 \wedge x_2) \vee (x_1 \wedge \bar{x}_2)) = (\bar{x}_1 \vee \bar{x}_2)$$

$$M_{13}(x_1, x_2) = ((\bar{x}_1 \wedge \bar{x}_2) \vee (\bar{x}_1 \wedge x_2) \vee (x_1 \wedge x_2)) = (\bar{x}_1 \vee x_2)$$

$$M_{14}(x_1, x_2) = ((\bar{x}_1 \wedge \bar{x}_2) \vee (x_1 \wedge \bar{x}_2) \vee (x_1 \wedge x_2)) = (x_1 \vee \bar{x}_2)$$

$$M_{15}(x_1, x_2) = ((\bar{x}_1 \wedge x_2) \vee (x_1 \wedge \bar{x}_2) \vee (x_1 \wedge x_2)) = (x_1 \vee x_2)$$

$$M_{16}(x_1, x_2) = ((\bar{x}_1 \wedge \bar{x}_2) \vee (\bar{x}_1 \wedge x_2) \vee (x_1 \wedge \bar{x}_2) \vee (x_1 \wedge x_2)) = 1$$

(b) The truth talbe for the funcitons in part (a) are

x_1	x_2	M_1	M_2	M_3	M_4	M_5	M_6	M_7	M_8	M_9	M_{10}	M_{11}	M_{12}	M_{13}	M_{14}	M_{15}	M_{16}
0	0	0	1	0	0	0	1	1	1	0	0	0	1	1	1	0	1
0	1	0	0	1	0	0	1	0	0	1	1	0	1	1	0	1	1
1	0	0	0	0	1	0	0	1	0	1	0	1	1	0	1	1	1
1	1	0	0	0	0	1	0	0	1	0	1	1	0	1	1	1	1

(c) $f_1(x_1, x_2) = M_{15}(x_1, x_2)$

$$f_2(x_1, x_2) = M_{12}(x_1, x_2)$$

$$f_3(x_1, x_2) = M_1(x_1, x_2)$$

$$f_4(x_1, x_2) = M_{16}(x_1, x_2)$$

3. (a) The number of elements in the domain of f is $16 = 4^2 = |B|^2$ (b) With two variables, there are $4^3 = 256$ different Boolean functions. With three variables, there are $4^8 = 65536$ different Boolean functions.

$$(c) \quad f(x_1, x_2) = (1 \wedge \bar{x}_1 \wedge \bar{x}_2) \vee (1 \wedge \bar{x}_1 \wedge x_2) \vee (1 \wedge x_1 \wedge \bar{x}_2) \vee (0 \wedge x_1 \wedge x_2)$$

(d) Consider $f: B^2 \rightarrow B$, defined by $f(0, 0) = 0$, $f(0, 1) = 1$, $f(1, 0) = a$, $f(1, 1) = a$, and $f(0, a) = b$, with the images of all other pairs in B^2 defined arbitrarily. This function is not a Boolean function. If we assume that it is Boolean function then f can be computed with a Boolean expression $M(x_1, x_2)$. This expression can be put into minterm normal form:

$$M(x_1, x_2) = (c_1 \wedge \bar{x}_1 \wedge \bar{x}_2) \vee (c_2 \wedge \bar{x}_1 \wedge x_2) \vee (c_3 \wedge x_1 \wedge \bar{x}_2) \vee (c_4 \wedge x_1 \wedge x_2)$$

$$f(0, 0) = 0 \Rightarrow M(0, 0) = 0 \Rightarrow c_1 = 0$$

$$f(0, 1) = 1 \Rightarrow M(0, 0) = 1 \Rightarrow c_1 = 1$$

$$f(1, 0) = a \Rightarrow M(0, 0) = a \Rightarrow c_1 = a$$

$$f(1, 1) = a \Rightarrow M(0, 0) = a \Rightarrow c_1 = a$$

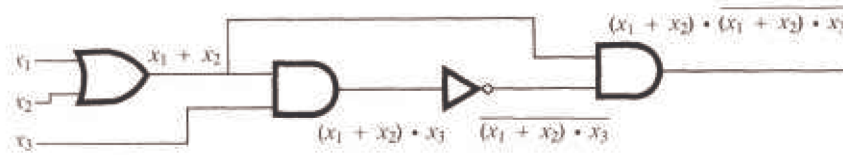
Therefore,

$$M(x_1, x_2) = (\bar{x}_1 \wedge x_2) \vee (a \wedge x_1 \wedge \bar{x}_2) \vee (a \wedge x_1 \wedge x_2)$$

$$M(0, a) = (\bar{0} \wedge a) \vee (a \wedge 0 \wedge \bar{a}) \vee (a \wedge 0 \wedge a) = a$$

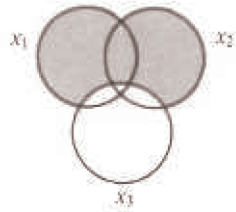
This contradicts $f(0, a) = b$, and so f is not a Boolean function.**Section 13.7**

1. (a)



$$\begin{aligned}
 (b) \quad f(x_1, x_2, x_3) &= \overline{((x_1 + x_2) \cdot x_3)} \cdot (x_1 + x_2) \\
 &= \overline{((x_1 + x_2) + \overline{x_3})} \cdot (x_1 + x_2) \\
 &= \overline{(x_1 + x_2)} \cdot (x_1 + x_2) + \overline{x_3} \cdot (x_1 + x_2) \\
 &= 0 + \overline{x_3} \cdot (x_1 + x_2) \\
 &= \overline{x_3} \cdot (x_1 + x_2)
 \end{aligned}$$

(c) The Venn diagram for the function is:



We can read off the minterm normal form from this diagram:

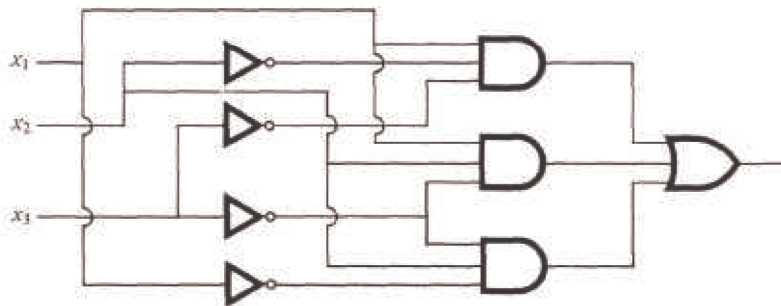
$$f(x_1, x_2, x_3) = x_1 \cdot \overline{x_2} \cdot \overline{x_3} + x_1 \cdot x_2 \cdot \overline{x_3} + \overline{x_1} \cdot x_2 \cdot \overline{x_3}$$

(d)

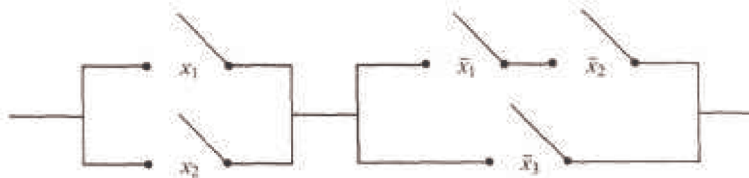
Simplified form:



Minterm form:



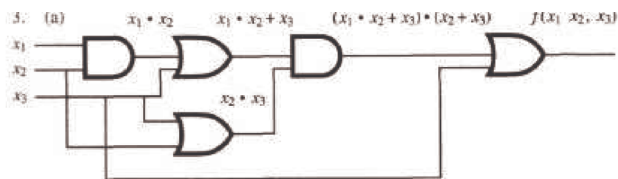
(e)



(f)

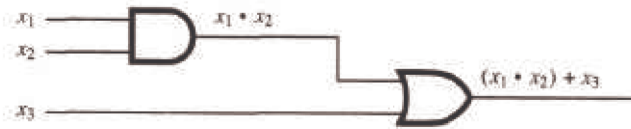
x_1	x_2	x_3	$x_1 + x_2$	$\overline{(x_1 + x_2)} \cdot x_3$	f
0	0	0	0	1	0
0	0	1	0	1	0
0	1	0	1	1	1
0	1	1	1	0	0
1	0	0	1	1	1
1	0	1	1	0	0
1	1	0	1	1	1
1	1	1	1	0	0

Current will flow only when one of the switches x_1 or x_2 is On and x_3 is Off.

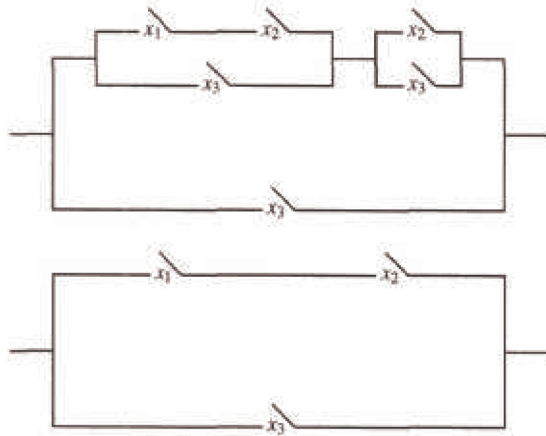


(b) $f(x_1, x_2, x_3) = (((x_1 \cdot x_2) + x_3) \cdot (x_2 + x_3)) + x_3$
 placing ()'s to indicate order of evaluation
 $= (((x_1 \cdot x_2) \cdot (x_2)) + x_3) + x_3$
 by the distributive law of $+$ over \cdot
 $= (x_1 \cdot (x_2 \cdot x_2)) + (x_3 + x_3)$
 by the associative laws of \cdot and $+$
 $= (x_1 \cdot x_2) + x_3$
 by the idempotent laws of \cdot and $+$

(c)

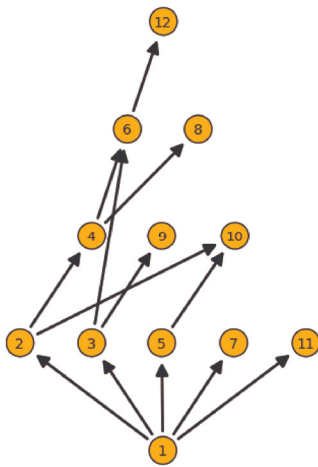


(d)

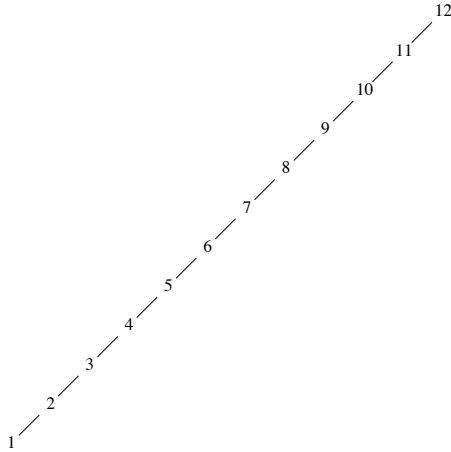
**Supplementary Exercises—Chapter 13**

1. (a) The following Sage input generates an ordering diagram.

```
Poset({1:[2,3,5,7,11],2:[4,6,10],3:[6,9],4:[6,8,12],5:[10],6:[12]}).plot()
```



(b) The ordering diagram for \leq is a chain



3. (a) $4 \vee 8 = 8, 3 \vee 15 = 15, 4 \wedge 8 = 4, 3 \wedge 15 = 3, 3 \wedge 5 = 3$.

(b) Yes. Let $a, b, c \in P$ and assume that there are n primes, p_1, p_2, \dots, p_n that appear as factors of a, b and c . Then we can write

$$a = p_1^{i_1} p_2^{i_2} \cdots p_n^{i_n}$$

$$b = p_1^{j_1} p_2^{j_2} \cdots p_n^{j_n}$$

$$c = p_1^{k_1} p_2^{k_2} \cdots p_n^{k_n}$$

where each exponent is a nonnegative integer. The greatest common divisor and least common multiple of two integers such as a and b can be expressed in terms of these exponents.

$$a \wedge b = \gcd(a, b) = p_1^{m_1} p_2^{m_2} \cdots p_n^{m_n}$$

where $m_r = \min(i_r, j_r)$ and

$$a \vee b = \text{lcm}(a, b) = p_1^{M_1} p_2^{M_2} \cdots p_n^{M_n}$$

where $M_r = \max(i_r, j_r)$.

Based on this observation, we can compare $a \wedge (b \vee c)$ and $(a \wedge b) \vee (a \wedge c)$. The exponent of p , is $\min(i_r, \max(j_r, k_r))$ in $a \wedge (b \vee c)$ and $\max(\min(i_r, j_r), \min(i_r, k_r))$ in $(a \wedge b) \vee (a \wedge c)$. These two exponents are equal; this is easiest to verify by checking the possible relative sizes of i_r, j_r and k_r . Therefore, the lattice is distributive.

(c) The least element is 1. There is no greatest element.

5. (a) The ordering diagram is the one-cube in Figure 9.4.5. It is interesting to note that the poset relation is really the logical implication, \Rightarrow , since $0 \Rightarrow 0, 0 \Rightarrow 1, 1 \Rightarrow 1$ are all true statements.

(b) From the definitions of lub and gcb and part (a) we have the tables

\wedge	0	1
0	0	0
1	0	1

\vee	0	1
0	0	1
1	1	1

which are the logical tables for the connectives "and" and "or."

(c) $L^2 = L \times L = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$ where the poset relation \leq on L^2 and the binary operations \wedge and \vee are all defined componentwise so that, for example, $(0, 1) \leq (1, 1)$, since in the two first coordinates, $0 \leq 1$ and in the two second coordinates, $1 \leq 1$. Also, for example, $(0, 1) \wedge (1, 0) = (0 \wedge 1, 1 \wedge 0) = (0, 0)$. The operation tables are given in the solution of Exercise 1 Section 13.5. The Hasse diagram for L^2 is the two-cube.

(d) The Hasse diagram for L^3 is the three-cube. Tables for \wedge and \vee can easily be constructed where, for example,

$$(1, 0, 0) \vee (0, 1, 0) = (1 \vee 0, 0 \vee 1, 0 \vee 0) = (1, 1, 0)$$

7. (a) No. It is not true that every pair of elements in A has both a *lub* and a *gib*

in A . For example, $10 \vee 4$ does not exist in A .

(b) Yes. For all $a, b \in A, a \neq b$,

$a \vee b =$ the maximum of a and b ,

$a \wedge b =$ the minimum of a and b .

9. $(x + y) \cdot (x + \bar{y}) = x + (y \cdot \bar{y})$ by the distributive law of $+$ over \cdot
 $= x + 0$ by the complement law
 $= x$ by the identity law

The switching circuit diagram has a single switch labeled x .

11. (a)

x	complement(s) of x
0	1
a_1	a_2, a_3, a_4, a_6
a_2	a_1, a_5
a_3	a_1, a_5
a_4	a_1, a_5
a_5	a_2, a_3, a_4, a_6
a_6	a_1, a_5
1	0

(b) No, it is not distributive, for if it were, complements would be unique.

13. (a) $D_{20} = \{1, 2, 4, 5, 10, 20\}$ contains 6 elements and so cannot be a Boolean algebra by Corollary 13.4.1.

(b) $D_{27} = \{1, 3, 9, 27\}$ has four elements and so we cannot use Corollary 13.4.1 to rule it out as a Boolean algebra. However, 3 has no complement, which means that D_{27} is not a Boolean algebra.

(c) $D_{35} = \{1, 5, 7, 35\}$ has $4 = 2^2$ elements, and so that it may be a Boolean algebra by Corollary 13.4.1. We can confirm through the definition of a Boolean algebra that it is.

(d) Notice that $210 = 2 \cdot 3 \cdot 5 \cdot 7$, which means that $|D_{210}| = 16 = 2^4$ and so Corollary 13.4.1 can't be used to rule it out as a Boolean algebra. Indeed, D_{210} is a Boolean algebra, which can be confirmed by applying the definition of a Boolean algebra.

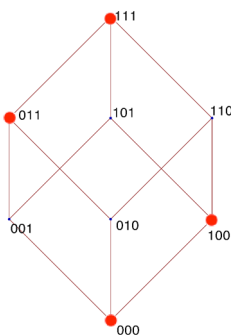
15. (a) First, by definition of subsystem in Section 11.5, a sub-Boolean algebra of a Boolean algebra B is a subset W of B which is a Boolean algebra under the same operations as B . Specifically, W must satisfy the conditions:

(i) The 0 and 1 of B must be in W ,

(ii) $a \in W \Rightarrow \bar{a} \in W$

(iii) $a, b \in W \Rightarrow a \vee b \in W$ and $a \wedge b \in W$.

Hence if W is to contain 4 elements it must be of the form $\{0, \beta, \bar{\beta}, 1\}$. $W_1 = \{(0, 0, 0), (0, 1, 1), (1, 0, 0), (1, 1, 1)\}$ is one such set. The 3-cube below illustrates this sub-Boolean algebra.



There are two others that are isomorphic to this one, where Corollary 13.4.2, assures us of this isomorphism.

(b) Again, the form of the sub-Boolean algebra with four elements must be $\{0, \beta, \bar{\beta}, 1\}$. Since the 2^n elements of B_2^n can be paired up with their complements to give us 2^{n-1} pairs, there are $2^{n-1} - 1$ ways to select the elements β and $\bar{\beta}$ (0 and its complement, 1, are already selected). Of course, all of these sub-Boolean algebras are isomorphic.

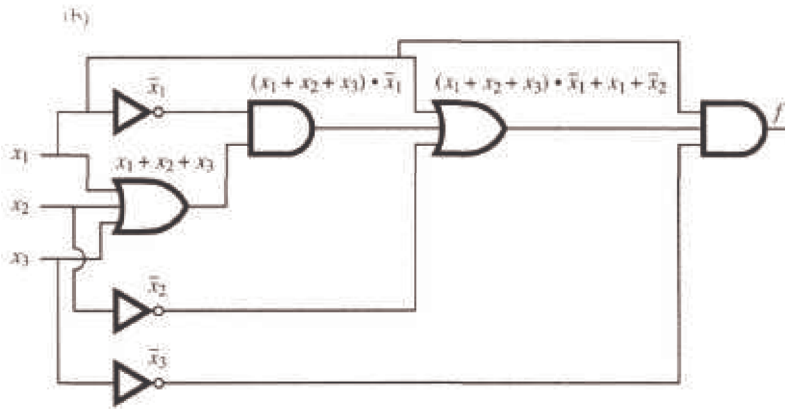
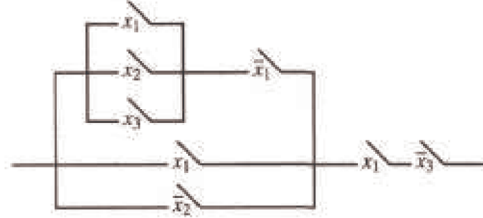
(c) A sub-Boolean algebra with 2^k elements must have k atoms; so the selection of k elements that will act as atoms can be considered in counting numbers of sub-Boolean algebras of a certain size. What is the number? We leave it to the reader in the general case.

17. $(\bar{x}_1 \wedge x_2 \wedge x_3) \vee (\bar{x}_1 \wedge \bar{x}_2 \wedge x_3) \vee (x_1 \wedge x_2 \wedge x_3)$

19. (a) Since each of the three variables can be any one of two values there are 2^3 rows, (See Table 13.6.3 for an example.) For n variables there are 2^n rows.

(b) For each row, there can be any one of two truth values. Since there are $2^3 = 8$ rows there are $2^8 = 256$ functions. For n variables and $m = 2^n$ rows, there are $2^m = 2^{2^n}$ functions.

21. (a)



$$\begin{aligned}
 (c) \quad f(x_1, x_2, x_3) &= ((x_1 + x_2 + x_3) \cdot \bar{x}_1 + x_1 + \bar{x}_2) \cdot x_1 \cdot \bar{x}_3 \\
 &= (x_1 \cdot \bar{x}_1 + x_2 \cdot \bar{x}_1 + x_3 \cdot \bar{x}_1 + x_1 + \bar{x}_2) \cdot x_1 \cdot \bar{x}_3 \\
 &= (0 + x_2 \cdot \bar{x}_1 + x_3 \cdot \bar{x}_1 + x_1 + \bar{x}_2) \cdot x_1 \cdot \bar{x}_3 \\
 &= (x_2 \cdot \bar{x}_1 + x_3 \cdot \bar{x}_1 + x_1 + \bar{x}_2) \cdot x_1 \cdot \bar{x}_3 \\
 &= x_2 \cdot \bar{x}_1 \cdot x_1 \cdot \bar{x}_3 + x_3 \cdot \bar{x}_1 \cdot x_1 \cdot \bar{x}_3 + x_1 \cdot x_1 \cdot \bar{x}_3 + \bar{x}_2 \cdot x_1 \cdot \bar{x}_3 \\
 &= x_2 \cdot 0 \cdot \bar{x}_3 + x_3 \cdot 0 \cdot \bar{x}_3 + x_1 \cdot \bar{x}_3 + \bar{x}_2 \cdot x_1 \cdot \bar{x}_3 \\
 &= x_1 \cdot \bar{x}_3 + \bar{x}_2 \cdot x_1 \cdot \bar{x}_3 \\
 &= x_1 \cdot \bar{x}_3 \cdot (1 + \bar{x}_2)
 \end{aligned}$$

Switching and gate diagrams to be added.

23. (a) $z = (\bar{x}_1 + x_2) + \bar{x}_2 \cdot \bar{x}_3$

(b) $z = (\bar{x}_1 + x_2) + \bar{x}_2 \cdot \bar{x}_3$
 $= (\bar{x}_1 + x_2) + (\bar{x}_2 + \bar{x}_3)$
 $= \bar{x}_1 + (x_2 + \bar{x}_2) + \bar{x}_3$
 $= \bar{x}_1 + 1 + \bar{x}_3$
 $= 1$

The circuit is always on, no gates are necessary.

CHAPTER 14

Section 14.1

1. (a) S_1 is not a submonoid since the identity of $[\mathbb{Z}_8, \times_8]$, which is 1, is not in S_1 . S_2 is a submonoid since $1 \in S_2$ and S_2 is closed under multiplication; that is, for all $a, b \in S_2$, $a \times_8 b$ is in S_2 .

(b) The identity of $\mathbb{N}^{\mathbb{N}}$ is the identity function $i: \mathbb{N} \rightarrow \mathbb{N}$ defined by $i(a) = a, \forall a \in \mathbb{N}$. If $a \in \mathbb{N}$, $i(a) = a \leq a$, thus the identity of $\mathbb{N}^{\mathbb{N}}$ is in S_1 . However, the image of 1 under any function in S_2 is 2, and thus the identity of $\mathbb{N}^{\mathbb{N}}$ is not in S_2 , so S_2 is not a submonoid. The composition of any two functions in S_1 , f and g , will be a function in S_1 :

$$(f \circ g)(n) = f(g(n)) \leq g(n) \text{ since } f \text{ is in } S_1 \\ \leq n \text{ since } g \text{ is in } S_1$$

Thus $f \circ g \in S_1$, and the two conditions of a submonoid are satisfied and S_1 is a submonoid of $\mathbb{N}^{\mathbb{N}}$.

(c) The first set is a submonoid, but the second is not since the null set has a non-finite complement.

3. The set of $n \times n$ real matrices is a monoid under matrix multiplication. This follows from the laws of matrix algebra in Chapter 5. To prove that the set of stochastic matrices is a monoid over matrix multiplication, we need only show that the identity matrix is stochastic (this is obvious) and that the set of stochastic matrices is closed under matrix multiplication. Let A and B be $n \times n$ stochastic matrices.

$$(AB)_{ij} = \sum_{k=1}^n a_{ik} b_{kj}$$

The sum of the j^{th} column is

$$\begin{aligned} \sum_{j=1}^n (AB)_{ij} &= \sum_{k=1}^n a_{1k} b_{kj} + \sum_{k=1}^n a_{1k} b_{kj} + \cdots + \sum_{k=1}^n a_{nk} b_{kj} \\ &= \sum_{k=1}^n (a_{1k} b_{kj} + a_{1k} b_{kj} + \cdots + a_{nk} b_{kj}) \\ &= \sum_{k=1}^n b_{kj} (a_{1k} + a_{1k} + \cdots + a_{nk}) \\ &= \sum_{k=1}^n b_{kj} \quad \text{since } A \text{ is stochastic} \\ &= 1 \quad \text{since } B \text{ is stochastic} \end{aligned}$$

Section 14.2

1. (a) For a character set of 350 symbols, the number of bits needed for each character is the smallest n such that 2^n is greater than or equal to 350. Since $2^9 = 512 > 350 > 2^8$, 9 bits are needed,

(b) $2^{12} = 4096 > 3500 > 2^{11}$; therefore, 12 bits are needed.

3. This grammar defines the set of all strings over B for which each string is a palindrome (same string if read forward or backward).

5. (a) Terminal symbols: The null string, 0, and 1.

Nonterminal symbols: S, E .

Starting symbol: S .

Production rules: $S \rightarrow 00S, S \rightarrow 01S, S \rightarrow 10S, S \rightarrow 11S, S \rightarrow E, E \rightarrow 0, E \rightarrow 1$

This is a regular grammar.

(b) Terminal symbols: The null string, 0, and 1.

Nonterminal symbols: S, A, B, C

Starting symbol: S

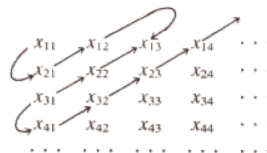
Production rules: $S \rightarrow 0A, S \rightarrow 1A, S \rightarrow \lambda, A \rightarrow 0B, A \rightarrow 1B, A \rightarrow \lambda, B \rightarrow 0C, B \rightarrow 1C, B \rightarrow A, C \rightarrow 0, C \rightarrow 1, C \rightarrow \lambda$

This is a regular grammar.

(c) See Exercise 3. This language is not regular.

7. If s is in A^* and L is recursive, we can answer the question "Is s in L^c ?" by negating the answer to "Is s in L ?"

9. (a) List the elements of each set x_i in a sequence $x_{i1}, x_{i2}, x_{i3}, \dots$

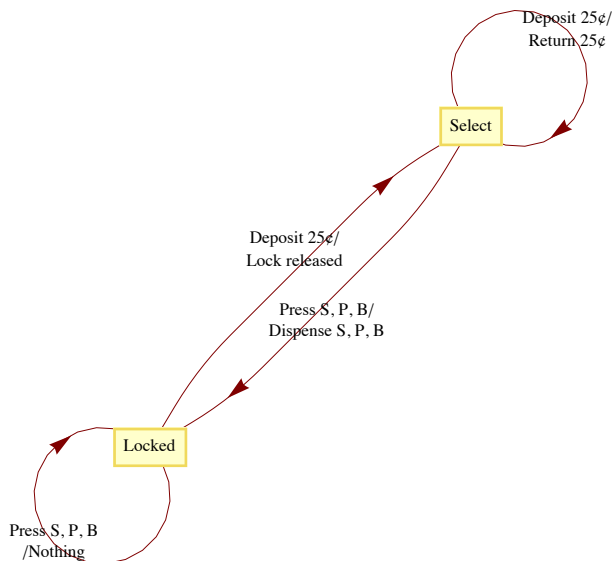


Then draw arrows as shown above and list the elements of the union in order established by this pattern: $x_{11}, x_{21}, x_{12}, x_{13}, x_{22}, x_{31}, x_{41}, x_{32}, x_{23}, x_{14}, x_{15}, \dots$

(b) Each of the sets A^1, A^2, A^3, \dots are countable and A^* is the union of these sets; hence A^* is countable.

Section 14.3

x	s	$Z(x, s)$	$t(x, s)$
Deposit 25 ¢	Locked	Nothing	Select
Deposit 25 ¢	Select	Return 25 ¢	Select
Press S	Locked	Nothing	Locked
Press S	Select	Dispense S	Locked
Press P	Locked	Nothing	Locked
Press P	Select	Dispense P	Locked
Press B	Locked	Nothing	Locked
Press B	Select	Dispense B	Locked



3. {000, 011, 101, 110, 111}

5. (a) Input: 10110, Output: 11011 \Rightarrow 10110 is in position 27

Input: 00100, Output: 00111 \Rightarrow 00100 is in position 7

Input: 11111, Output: 10101 \Rightarrow 11111 is in position 21

(b) Let $x = x_1 x_2 \dots x_n$ and recall that for $n \geq 1$, $G_{n+1} = \begin{pmatrix} 0 & G_n \\ 1 & G_n^r \end{pmatrix}$, where G_n^r is the reverse of G_n . To prove that the Gray Code Decoder always works, let $p(n)$ be the proposition "Starting in Copy state, x 's output is the position of x in G_n ; and starting in Complement state, x 's output is the position of x in G_n^r ." That $p(1)$ is true is easy to verify for both possible values of x , 0 and 1. Now assume that for some $n \geq 1$, $p(n)$ is true and consider $x = x_1 x_2 \dots x_n x_{n+1}$.

If $x_1 = 0$,

x 's output = 0 followed by $(x_2 \dots x_n x_{n+1})$'s output starting in Copy
 = 0 followed by $(x_2 \dots x_n x_{n+1})$'s position in G_n
 = x 's position in G_{n+1}

If $x_1 = 1$,

x 's output = 1 followed by $(x_2 \dots x_n x_{n+1})$'s output starting in Complement
 = 1 followed by $(x_2 \dots x_n x_{n+1})$'s position in G_n^r
 = x 's position in G_{n+1}

Section 14.4

Input String	a	b	c	aa	ab	ac
1.						
1	$(a, 1)$	$(a, 2)$	$(c, 3)$	$(a, 1)$	$(a, 2)$	$(c, 3)$
2	$(a, 2)$	$(a, 1)$	$(c, 3)$	$(a, 2)$	$(a, 1)$	$(c, 3)$
3	$(c, 3)$	$(c, 3)$	$(c, 3)$	$(c, 3)$	$(c, 3)$	$(c, 3)$
Input String	ba	bb	bc	ca	cb	cc
1	$(a, 2)$	$(a, 1)$	$(c, 3)$	$(c, 3)$	$(c, 3)$	$(c, 3)$
2	$(a, 1)$	$(a, 2)$	$(c, 3)$	$(c, 3)$	$(c, 3)$	$(c, 3)$
3	$(c, 3)$	$(c, 3)$	$(c, 3)$	$(c, 3)$	$(c, 3)$	$(c, 3)$

We can see that $T_a T_a = T_{aa} = T_a$, $T_a T_b = T_{ab} = T_b$, etc. Therefore, we have the following monoid:

	T_a	T_b	T_b
T_a	T_a	T_b	T_c
T_b	T_b	T_a	T_c
T_c	T_c	T_c	T_c

Notice that T_a is the identity of this monoid.

	Input String	1	2	11	12	21	22		
(b)	A	C	B	A	D	D	A		
	B	D	A	B	C	C	B		
	C	A	D	C	B	B	C		
	D	B	C	D	A	A	D		
	Input String	111	112	121	122	211	212	221	222
	A	C	B	B	C	B	C	C	B
	B	D	A	A	D	A	D	D	A
	C	B	C	C	B	C	B	B	C
	D	B	C	C	B	C	B	B	C

We have the following monoid:

	T_1	T_2	T_{11}	T_{12}
T_1	T_{11}	T_{12}	T_1	T_2
T_2	T_b	T_{11}	T_2	T_1
T_{11}	T_1	T_2	T_{11}	T_{12}
T_{12}	T_2	T_1	T_{12}	T_{11}

Notice that T_{11} is the identity of this monoid.

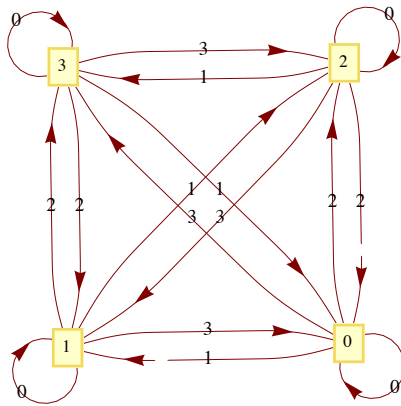
3. Yes, just consider the unit time delay machine of Figure 14.4.2. Its monoid is described by the table at the end of Section 14.4 where the T_λ row and T_λ column are omitted. Next consider the machine in Figure 14.5.3. The monoid of this machine is:

	T_λ	T_0	T_1	T_{00}	T_{01}	T_{10}	T_{11}
T_λ	T_λ	T_0	T_1	T_{00}	T_{01}	T_{10}	T_{11}
T_0	T_0	T_{00}	T_{01}	T_{00}	T_{01}	T_{10}	T_{11}
T_1	T_1	T_{10}	T_{11}	T_{00}	T_{01}	T_{10}	T_{11}
T_{00}	T_{00}	T_{00}	T_{01}	T_{00}	T_{01}	T_{10}	T_{11}
T_{01}	T_{01}	T_{10}	T_{11}	T_{00}	T_{01}	T_{10}	T_{11}
T_{10}	T_{10}	T_{00}	T_{01}	T_{00}	T_{01}	T_{10}	T_{11}
T_{11}	T_{11}	T_{10}	T_{11}	T_{00}	T_{01}	T_{10}	T_{11}

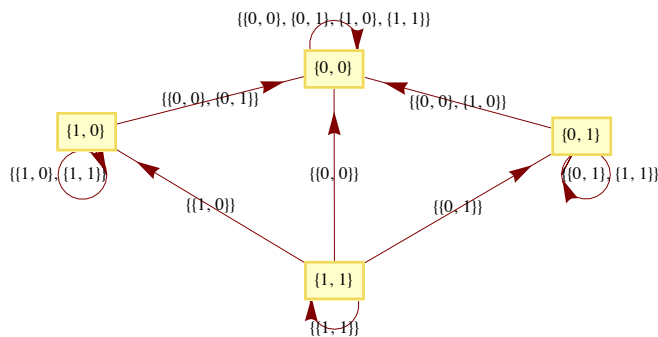
Hence both of these machines have the same monoid, however, their transition diagrams are nonisomorphic since the first has two vertices and the second has seven.

Section 14.5

1. (a)



(b)



Supplementary Exercises—Chapter 14

1. Let $f, g, h \in M$, and $a \in B$.

$$\begin{aligned}
 ((f * g) * h)(a) &= (f * g)(a) \wedge h(a) \\
 &= (f(a) \wedge g(a)) \wedge h(a) \\
 &= f(a) \wedge (g(a) \wedge h(a)) \\
 &= f(a) \wedge (g * h)(a) \\
 &= (f * (g * h))(a)
 \end{aligned}$$

Therefore $(f * g) * h = f * (g * h) \Rightarrow *$ is associative.

The identity for $*$ is the function $u \in M$ where $u(a) = 1$ = the “one” of B . If $a \in B$

$$(f * u)(a) = f(a) \wedge u(a) = f(a) \wedge 1 = f(a)$$

Therefore $f * u = f$. Similarly $u * f = f$.

There are $2^2 = 4$ functions in M for $B = B_2$. These four functions are named in the text (see Figure 14.1.1). The table for $*$ is

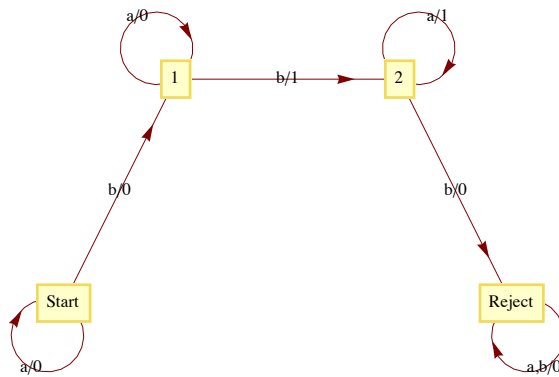
	z	i	t	u
z	z	z	z	z
i	z	i	z	i
t	z	z	t	t
u	z	u	t	u

3. $\{a, bb, bbb, bbbb, \dots\}$

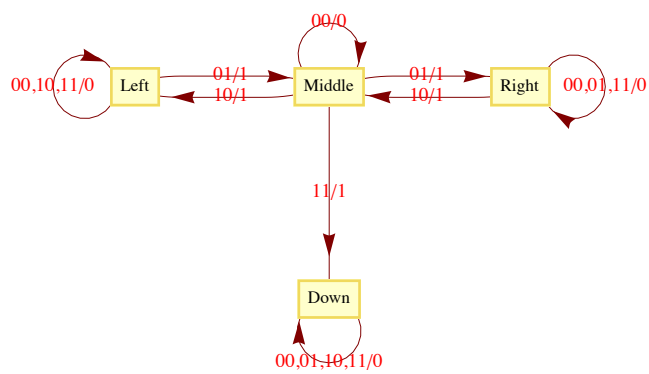
5. S = start symbol. Nonterminals = $\{S, B_0, B_1, B_2\}$

$$\begin{aligned}
 S &\rightarrow B_0 & B_0 &\rightarrow a B_0 & B_0 &\rightarrow b B_1 \\
 B_1 &\rightarrow a B_1 & B_1 &\rightarrow b B_2 & B_1 &\rightarrow b \\
 B_2 &\rightarrow a B_2 & B_2 &\rightarrow a
 \end{aligned}$$

7.

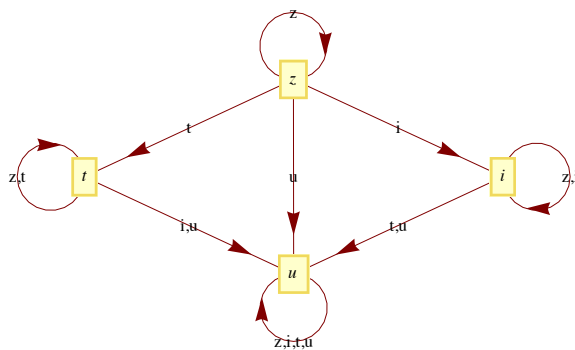


9. (a)



(b) The possible output sequences are 100, 010, 001, and 111. Note: Output for $t = 3$ is determined by the next state, $s(4)$. If $s(4) = s(3)$, output at $t = 3$ is 0, while if $s(4) \neq s(3)$, output at $t = 3$ is 1.

11.



CHAPTER 15

Section 15.1

1. The only other generator is -1 .

3. If $|G| = m$, $m > 2$, and $G = \langle a \rangle$, then $a, a^2, \dots, a^{m-1}, a^m = e$ are distinct elements of G . Furthermore, $a^{-1} = a^{m-1} \neq a$. If $1 \leq k \leq m$, a^{-1} generates a^k :

$$\begin{aligned}(a^{-1})^{m-k} &= (a^{m-1})^{m-k} = a^{m^2-m-k+k} \\ &= (a^m)^{m-k-1} * a^k = e * a^k = a^k\end{aligned}$$

Similarly, if G is infinite and $G = \langle a \rangle$, then a^{-1} generates G .

5. (a) No. Assume that $q \in \mathbb{Q}$ generates \mathbb{Q} . Then $\langle q \rangle = \{nq : n \in \mathbb{Z}\}$. But this gives us at most integer multiples of q , not every element in \mathbb{Q} .

(b) No. Similar reasoning to part a.

(c) Yes. 6 is a generator of $6\mathbb{Z}$.

(d) No.

(e) Yes, $(1, 1, 1)$ is a generator of the group.

7. Theorem 15.1.4 implies that a generates \mathbb{Z}_n if and only if the greatest common divisor of n and a is 1 (i. e., n and a are relatively prime). Therefore the list of generators of \mathbb{Z}_n are the integers in \mathbb{Z}_n that are relatively prime to n . The generators of \mathbb{Z}_{25} are all of the nonzero elements except 5, 10, 15, and 20. The generators of \mathbb{Z}_{256} are the odd integers in \mathbb{Z}_{256} since 256 is 2^8 . *Mathematica* expression to generate these sets are

```
Select[Range[0, 24], Function[a, GCD[25, a] == 1]]
```

```
{1, 2, 3, 4, 6, 7, 8, 9, 11, 12, 13, 14, 16, 17, 18, 19, 21, 22, 23, 24}
```

```
Select[Range[0, 255], Function[a, GCD[256, a] == 1]]
```

```
{1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25, 27, 29, 31, 33, 35, 37, 39, 41, 43, 45, 47, 49, 51, 53, 55, 57, 59, 61, 63, 65,
 67, 69, 71, 73, 75, 77, 79, 81, 83, 85, 87, 89, 91, 93, 95, 97, 99, 101, 103, 105, 107, 109, 111, 113, 115, 117, 119,
 121, 123, 125, 127, 129, 131, 133, 135, 137, 139, 141, 143, 145, 147, 149, 151, 153, 155, 157, 159, 161, 163, 165,
 167, 169, 171, 173, 175, 177, 179, 181, 183, 185, 187, 189, 191, 193, 195, 197, 199, 201, 203, 205, 207, 209, 211,
 213, 215, 217, 219, 221, 223, 225, 227, 229, 231, 233, 235, 237, 239, 241, 243, 245, 247, 249, 251, 253, 255}
```

9. (a) $\theta: \mathbb{Z}_{77} \rightarrow \mathbb{Z}_7 \times \mathbb{Z}_{11}$

$$21 \rightarrow (0, 10)$$

$$5 \rightarrow (5, 5)$$

$$7 \rightarrow (0, 7)$$

$$15 \rightarrow (1, 4)$$

$$\text{sum} = 48 \leftarrow (6, 4) = \text{sum}$$

The final sum, 48, is obtained by using the facts that $\theta^{-1}(1, 0) = 22$ and $\theta^{-1}(0, 1) = 56$

$$\begin{aligned}\theta^{-1}(6, 4) &= 6 \times_{77} \theta^{-1}(1, 0) + 4 \times_{77} \theta^{-1}(0, 1) \\ &= 6 \times_{77} 22 +_{77} 4 \times_{77} 56 \\ &= 55 +_{77} 70 \\ &= 48\end{aligned}$$

(b) Using the same isomorphism:

$$25 \rightarrow (4, 3)$$

$$26 \rightarrow (5, 4)$$

$$40 \rightarrow (5, 7)$$

$$\text{sum} = (0, 3)$$

$$\begin{aligned}\theta^{-1}(0, 3) &= 3 \times_{77} \theta^{-1}(0, 1) \\ &= 3 \times_{77} 56 \\ &= 14\end{aligned}$$

The actual sum is 91. Our result is incorrect, since 91 is not in \mathbb{Z}_{77} . Notice that 91 and 14 differ by 77. Any error that we get using this technique will be a multiple of 77.

Section 15.2

1. Call the subsets A and B respectively. If we choose $0 \in A$ and $5 \in B$ we get $0 +_{10} 5 = 5 \in B$. On the other hand, if we choose $3 \in A$ and $8 \in B$, we get $3 +_{10} 8 = 1 \in A$. Therefore, the induced operation is not well defined on $\{A, B\}$.

3. (a) The four distinct cosets in G/H are

$$H = \{(0, 0), (2, 0)\}$$

$$(1, 0) + H = \{(1, 0), (3, 0)\}$$

$$(0, 1) + H = \{(0, 1), (2, 1)\},$$

$$\text{and } (1, 1) + H = \{(1, 1), (3, 1)\}$$

None of these cosets generates G/H ; therefore G/H is not cyclic. Hence G/H must be isomorphic to $\mathbb{Z}_2 \times \mathbb{Z}_2$.

(b) The factor group is isomorphic to $[\mathbb{R}; +]$. Each coset of \mathbb{R} is a line in the complex plane that is parallel to the x-axis: $\tau : \mathbb{C}/\mathbb{R} \rightarrow \mathbb{R}$, where $T(\{a + bi \mid a \in \mathbb{R}\}) = b$ is an isomorphism.

(c) $\langle 8 \rangle = \{0, 4, 8, 12, 16\} \Rightarrow |\mathbb{Z}_{20}/\langle 8 \rangle| = 4$.

The four cosets are: $\bar{0}, \bar{1}, \bar{2}$, and $\bar{3}$. $\bar{1}$ generates all four cosets. The factor group is isomorphic to $[\mathbb{Z}_4, +_4]$ because $\bar{1}$ generates it.

$$\begin{aligned} 5. \quad a \in bH &\Leftrightarrow a = b * h \text{ for some } h \in H \\ &\Leftrightarrow b^{-1} * a = h \text{ for some } h \in H \\ &\Leftrightarrow b^{-1} * a \in H \end{aligned}$$

Section 15.3

$$1. \quad (a) \begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 4 & 3 & 2 \end{pmatrix} \quad (b) \begin{pmatrix} 1 & 2 & 3 & 4 \\ 4 & 3 & 1 & 2 \end{pmatrix}$$

$$(c) \begin{pmatrix} 1 & 2 & 3 & 4 \\ 3 & 4 & 2 & 1 \end{pmatrix} \quad (d) \begin{pmatrix} 1 & 2 & 3 & 4 \\ 3 & 4 & 2 & 1 \end{pmatrix}$$

$$(e) \begin{pmatrix} 1 & 2 & 3 & 4 \\ 4 & 2 & 1 & 3 \end{pmatrix} \quad (f) \begin{pmatrix} 1 & 2 & 3 & 4 \\ 3 & 1 & 4 & 2 \end{pmatrix}$$

$$(g) \begin{pmatrix} 1 & 2 & 3 & 4 \\ 2 & 1 & 4 & 3 \end{pmatrix}$$

3. Yes and no, respectively

$$5. D_4 = \{i, r, r^2, r^3, f_1, f_2, f_3, f_4\}$$

Where i is the identity function, $r = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 2 & 3 & 4 & 1 \end{pmatrix}$, and

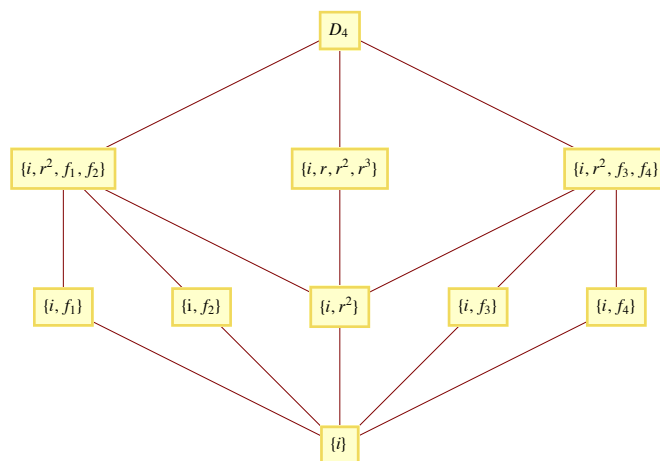
$$f_1 = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 4 & 3 & 2 & 1 \end{pmatrix} \quad f_2 = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 2 & 1 & 4 & 3 \end{pmatrix}$$

$$f_3 = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 3 & 2 & 1 & 4 \end{pmatrix} \quad f_4 = \begin{pmatrix} 1 & 2 & 3 & 4 \\ 1 & 4 & 3 & 2 \end{pmatrix}$$

The operation table for the group is

\circ	i	r	r^2	r^3	f_1	f_2	f_3	f_4
i	i	r	r^2	r^3	f_1	f_2	f_3	f_4
r	r	r^2	r^3	i	f_4	f_3	f_1	f_2
r^2	r^2	r^3	i	r	f_2	f_1	f_4	f_3
r^3	r^3	i	r	r^2	f_3	f_4	f_2	f_1
f_1	f_1	f_3	f_2	f_4	i	r^2	r	r^3
f_2	f_2	f_4	f_1	f_3	r^2	i	r^3	r
f_3	f_3	f_2	f_4	f_1	r^3	r	i	r^2
f_4	f_4	f_1	f_3	f_2	r	r^3	r^2	i

A lattice diagram of its subgroups is



All proper subgroups are cyclic except $\{i, r^2, f_1, f_2\}$ and $\{i, r^2, f_3, f_4\}$. Each 2-element subgroup is isomorphic to \mathbb{Z}_2 ; $\{i, r, r^2, r^3\}$ is isomorphic to \mathbb{Z}_4 ; and $\{i, r^2, f_1, f_2\}$ and $\{i, r^2, f_3, f_4\}$ are isomorphic to $\mathbb{Z}_2 \times \mathbb{Z}_2$.

7. One solution is to cite Exercise 3 at the end of Section 11.3. It can be directly applied to this problem. An induction proof of the problem at hand would be almost identical to the proof of the more general statement.

$$\begin{aligned} (t_1 t_2 \cdots t_r)^{-1} &= t_r^{-1} \cdots t_2^{-1} t_1^{-1} && \text{by Exercises 3 of Section 11.3} \\ &= t_r \cdots t_2 t_1 && \text{since each transposition inverts itself. } \blacksquare \end{aligned}$$

9. Part I: That $|S_k| = k!$ follows from Exercise 3 of Section 7.3.

Part II: Let f be the function defined on $\{1, 2, \dots, n\}$ by $f(1) = 2$, $f(2) = 3$, $f(3) = 1$, and $f(j) = j$ for $4 \leq j \leq n$; and let g be defined by $g(1) = 1$, $g(2) = 3$, $g(3) = 2$, and $g(j) = j$ for $4 \leq j \leq n$. Note that f and g are elements of S_n . Next, $(f \circ g)(1) = f(g(1)) = f(1) = 2$, while $(g \circ f)(1) = g(f(1)) = g(2) = 3$, hence $f \circ g \neq g \circ f$ and S_n is non-abelian for any $n \geq 3$.

11. (a) Both groups are non-abelian and of order 6; so they must be isomorphic, since only one such group exists up to isomorphism. The function $\theta: S_3 \rightarrow R_3$ defined by

$$\begin{aligned} \theta(i) &= I & \theta(f_1) &= F_1 \\ \theta(r_1) &= R_1 & \theta(f_2) &= F_2 \\ \theta(r_2) &= R_2 & \theta(f_3) &= F_3 \end{aligned}$$

is an isomorphism,

(b) Recall that since every function is a relation, it is natural to translate functions to Boolean matrices. Suppose that $f \in S_n$. We will define its image, $\theta(f)$, by

$$\theta(f)_{kj} = 1 \iff f(j) = k$$

That θ is a bijection follows from the existence of θ^{-1} . If A is a rook matrix,

$$\begin{aligned} \theta^{-1}(A)(j) = k &\iff \text{The 1 in column } j \text{ of } A \text{ appears in row } k \\ &\iff A_{kj} = 1 \end{aligned}$$

For $f, g \in S_n$,

$$\begin{aligned} \theta(f \circ g)_{kj} = 1 &\iff (f \circ g)(j) = k \\ &\iff \exists l \text{ such that } g(j) = l \text{ and } f(l) = k \\ &\iff \exists l \text{ such that } \theta(g)_{lj} = 1 \text{ and } \theta(f)_{kl} = 1 \\ &\iff (\theta(f) \theta(g))_{kj} = 1 \end{aligned}$$

Therefore, θ is an isomorphism. \blacksquare

Section 15.4

1. (a) Yes, the kernel is $\{1, -1\}$

(b) No, since $\theta_2(2 +_5 4) = \theta_2(1) = 1$, but $\theta_2(2) +_2 \theta_2(4) = 0 +_2 0 = 0$

(c) Yes, the kernel is $\{(a, -a) \mid a \in \mathbb{R}\}$

(d) No

3. $\langle r \rangle = \{i, r, r^2, r^3\}$ is a normal subgroup of D_4 . To see you could use the table given in the solution of Exercise 5 of Section 15.3 and verify that $a^{-1} h a \in \langle r \rangle$ for all $a \in D_4$ and $h \in \langle r \rangle$. A more efficient approach is to prove the general theorem that if H is a subgroup G with exactly two distinct left cosets, then H is normal.

$\langle f_1 \rangle$ is not a normal subgroup of D_4 . $\langle f_1 \rangle = \{i, f_1\}$ and if we choose $a = r$ and $h = f_1$ then $a^{-1} h a = r^3 f_1 r = f_2 \notin \langle f_1 \rangle$

5. $(\beta \circ \alpha)(a_1, a_2, a_3) = 0$ and so $\beta \circ \alpha$ is the trivial homomorphism, but a homomorphism nevertheless.

7. Let $x, y \in G$.

$$\begin{aligned} q(x * y) &= (x * y)^2 \\ &= x * y * x * y \\ &= x * x * y * y \text{ since } G \text{ is abelian} \\ &= x^2 * y^2 \\ &= q(x) * q(y) \end{aligned}$$

Hence, q is a homomorphism.

In order for q to be an isomorphism, it must be the case that no element other than the identity is its own inverse.

$$\begin{aligned} x \in \text{Ker}(q) &\Leftrightarrow q(x) = e \\ &\Leftrightarrow x * x = e \\ &\Leftrightarrow x^{-1} = x \end{aligned}$$

9. Proof: Recall: The inverse image of H' under θ is

$$\theta^{-1}(H') = \{g \in G \mid \theta(g) \in H'\}$$

Closure: Let $g_1, g_2 \in \theta^{-1}(H')$, then $\theta(g_1), \theta(g_2) \in H'$. Since H' is a subgroup of G' ,

$$\theta(g_1) \diamond \theta(g_2) = \theta(g_1 * g_2) \Rightarrow g_1 * g_2 \in \theta^{-1}(H')$$

Identity: By Theorem 15.4.2(a), $e \in \theta^{-1}(H')$.

Inverse: Let $a \in \theta^{-1}(H')$. Then $\theta(a) \in H'$ and by Theorem 15.4.2(b), $\theta(a)^{-1} = \theta(a^{-1}) \in H'$ and so $a^{-1} \in \theta^{-1}(H')$.

Section 15.5

1. (a) Error detected, since an odd number of 1s was received; ask for retransmission.

(b) No error detected; accept this block.

(c) No error detected; accept this block.

3. (a) Syndrome = (1, 0, 1). Corrected message = (1, 1, 0).

(b) Syndrome = (1, 1, 0). Corrected message = (0, 0, 1).

(c) Syndrome (0, 0, 0). Corrected message = received message.
= (0, 1, 1)

(d) Syndrome = (1, 1, 0). Corrected message = (1, 0, 0).

(e) Syndrome = (1, 1, 1). This syndrome occurs only if two bits have been switched. No reliable correction is possible.

5. Let G be the 9×10 matrix obtained by augmenting the 9×9 identity matrix with a column of ones. The function $e : \mathbb{Z}_2^9 \rightarrow \mathbb{Z}_2^{10}$ defined by $e(a) = aG$ will allow us to detect single errors, since $e(a)$ will always have an even number of ones.

Supplementary Exercises—Chapter 15

1. Theorem 15.1.3 guarantees that all subgroups of any cyclic group can be determined by finding all cyclic subgroups. We can find all cyclic subgroups of noncyclic groups but there may be other subgroups.

3. First, write 120 as a product of powers of distinct primes: $120 = 2^3 \cdot 3 \cdot 5$. The Chinese Remainder Theorem states that $\theta : \mathbb{Z}_{120} \rightarrow \mathbb{Z}_8 \times \mathbb{Z}_3 \times \mathbb{Z}_5$ defined by $\theta(k) = (k \bmod 8, k \bmod 3, k \bmod 5)$ is an isomorphism. In particular, $\theta(74) = (2, 2, 4)$ and $\theta(85) = (5, 1, 0)$. Therefore,

$$\begin{aligned}\theta(74 +_{120} 85) &= \theta(74) + \theta(85) \\ &= (2, 2, 4) + (5, 1, 0) \\ &= (7, 0, 4)\end{aligned}$$

Since $\theta(105) = (1, 0, 0)$, and $\theta(96) = (0, 0, 1)$, we can compute

$$\begin{aligned}\theta^{-1}(7, 0, 4) &= 7 \times_{120} 105 +_{120} 4 \times_{120} 96 \\ &= 39\end{aligned}$$

$$5. H = 0 + H = \{0, 4, 8\} = 4 + H = 8 + H$$

$$1 + H = \{1, 5, 9\} = 5 + H = 9 + H$$

$$2 + H = \{2, 6, 10\} = 6 + H = 10 + H$$

$$3 + H = \{3, 7, 11\} = 7 + H = 11 + H$$

The operation table for this factor group is the same as that of $[\mathbb{Z}_4, +_4]$ with k replaced with $k + H$.

7. (a) $|\mathbb{Z}_8| = 8$ and $|\langle 2 \rangle| = 4$, therefore there are 2 distinct left cosets, and they are:

$$0 + \langle 2 \rangle = \{0, 2, 4, 6\} = 2 + \langle 2 \rangle = 4 + \langle 2 \rangle = 6 + \langle 2 \rangle$$

$$1 + \langle 2 \rangle = \{1, 3, 5, 7\} = 3 + \langle 2 \rangle = 5 + \langle 2 \rangle = 7 + \langle 2 \rangle$$

(b) $|\mathbb{Z}_{12}| = 12$ and $|\langle 2 \rangle| = 6$, therefore there are 2 distinct left cosets and they are:

$$0 + \langle 2 \rangle = \{0, 2, 4, 6, 8, 10\} = 2 + \langle 2 \rangle = 4 + \langle 2 \rangle = 6 + \langle 2 \rangle = 8 + \langle 2 \rangle = 10 + \langle 2 \rangle$$

$$\text{and } 1 + \langle 2 \rangle = \{1, 3, 5, 7, 9, 11\} = 3 + \langle 2 \rangle = 5 + \langle 2 \rangle = 7 + \langle 2 \rangle = 9 + \langle 2 \rangle = 11 + \langle 2 \rangle$$

(c) Since both groups are of order 2 and there is only one group of order 2 up to isomorphism, they are isomorphic. A simpler group is \mathbb{Z}_2 .

7. Assume f is even, $f = t_1 \circ t_2 \circ \cdots \circ t_{2r}$ for some r , where each t_i is a transposition. Hence

$$f^{-1} = (t_1 \circ t_2 \circ \cdots \circ t_{2r})^{-1} = t_{2r} \circ \cdots \circ t_2 \circ t_1 \text{ by Exercise 11 of Section 15.3.}$$

Since the alternative, that f is odd, leads to f^{-1} being odd, f is even if and only if f^{-1} is even.

11. (a) This following is the "standard definition" of a Boolean algebra homomorphism.

$f : B_1 \rightarrow B_2$ is a Boolean algebra homomorphism if and only if for all $a, b \in B_1$.

$$(1) f(a \wedge b) = f(a) \wedge f(b)$$

$$(2) f(a \vee b) = f(a) \vee f(b)$$

$$(3) f(\bar{a}) = \overline{f(a)}$$

$$\begin{aligned}(b) (i) f(0) &= f(a \wedge \bar{a}) \\ &= f(a) \wedge f(\bar{a}) \\ &= f(a) \wedge \overline{f(a)} \\ &= 0\end{aligned}$$

and

$$\begin{aligned}f(1) &= f(a \vee \bar{a}) \\ &= f(a) \vee f(\bar{a}) \\ &= f(a) \vee \overline{f(a)} \\ &= 1\end{aligned}$$

Note : The 0 and 1 of B_1 may be different than those of B_2 .

(ii) $a \leq b \Rightarrow a = a \wedge b$ by Supplementary Exercise 4 of Chapter 13

$$\Rightarrow f(a) = f(a \wedge b) = f(a) \wedge f(b)$$

$$\Rightarrow f(a) \leq f(b) \text{ by the same exercise cited above.}$$

(iii) See the solution to Exercise 15 of the Supplementary section of Chapter 13 for the definition of Boolean subalgebra. Part (i) of this exercise shows that $f(B_1)$ contains the 0 and 1 of B_2 . The definition in part a shows that $f(a) \in f(B_1)$ has a complement, namely $f(\bar{a}) \in f(B_1)$, and also that $f(B_1)$ must be closed with respect to both \wedge and \vee . For example, if $a, b \in B_1$, then $a \wedge b \in B_1$, and since $f(a) \wedge f(b) = f(a \wedge b)$, $f(a) \wedge f(b) \in f(B_1)$.

$$13 \text{ (a)} \begin{pmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{pmatrix}$$

$$(b) \ e(1111) = 1111111 \text{ and } e(1001) = 1001001$$

(c) (i) Syndrome = 101 \Rightarrow Error in second bit, since 101 is the second row of P .

Corrected message = 0000.

(ii) Syndrome = 000 \Rightarrow No error in transmission. Correct message is 1010.

(iii) Syndrome = 001 \Rightarrow Error in seventh bit, since 001 is the seventh row of P .

Corrected message = 1011. (Since the error was in a parity bit, the actual message is not corrected.)

(d) The most direct way of proving that all single errors can be corrected is to compute the syndromes of each of the seven possible one-bit errors. Since each of them produces a distinct syndrome (the rows of P), single errors can always be corrected.

CHAPTER 16

Section 16.1

1. All but rings c and e are commutative. All of the rings have a unity element. The number 1 is the unity for all of the rings except c and e . The unity for $M_{2 \times 2}(\mathbb{R})$ is the two by two identity matrix; the unity for $M_{n \times n}(\mathbb{R})$ is the n by n identity matrix. The units are as follows:

(a) $\{1, -1\}$

(b) \mathbb{C}^*

(c) $\{A \mid |A| = \pm 1\}$

(d) \mathbb{Q}^*

(e) $\{A \mid A_{11}A_{22} - A_{12}A_{21} \neq 0\}$

(f) $\{1\}$

3. Hints: (a) Consider commutativity

(b) Solve $x^2 = 3x$ in both rings.

5. (a) We already know that $3\mathbb{Z}$ is a subgroup of the group \mathbb{Z} ; so part 1 of Theorem 16.1.1 is satisfied. We need only show that part 2 of the theorem holds: Let $3m, 3n \in 3\mathbb{Z}$.

$$(3m)(3n) = 3(3mn) \in 3\mathbb{Z}, \text{ since } 3mn \in \mathbb{Z}. \blacksquare$$

(b) The proper subrings are $\{0, 2, 4, 6\}$ and $\{0, 4\}$; while $\{0\}$ and \mathbb{Z}_8 are improper subrings.

(c) The proper subrings are $\{00, 01\}$, $\{00, 10\}$, and $\{00, 11\}$; while $\{00\}$ and $\mathbb{Z}_2 \times \mathbb{Z}_2$ are improper subrings.

7. (a) The left-hand side of the equation factors into the product $(x-2)(x-3)$. Since \mathbb{Z} is an integral domain, $x=2$ and $x=3$ are the only possible solutions.

(b) Over \mathbb{Z}_{12} , 2, 3, 6, and 11 are solutions. Although the equation factors into $(x-2)(x-3)$, this product can be zero without making x either 2 or 3. For example. If $x=6$ we get $(6-2) \times_{12} (6-3) = 4 \times_{12} 3 = 0$. Notice that 4 and 3 are divisors of zero.

9. Let R_1, R_2 , and R_3 be any rings, then

(a) R_1 is isomorphic to R_1 and so “is isomorphic to” is a reflexive relation on rings,

(b) R_1 is isomorphic to $R_2 \Rightarrow R_2$ is isomorphic to R_1 , and so “is isomorphic to” is a symmetric relation on rings,

(c) R_1 is isomorphic to R_2 , and R_2 is isomorphic to R_3 implies that R_1 is isomorphic to R_3 , and so “is isomorphic to” is a transitive relation on rings.

We haven’t proven these properties here, just stated them. The combination of these observations implies that “is isomorphic to” is an equivalence relation on rings,

11. (a) Commutativity is clear from examination of a multiplication table for $\mathbb{Z}_2 \times \mathbb{Z}_3$. More generally, we could prove a theorem that the direct product of two or more commutative rings is commutative. $(1, 1)$ is the unity of $\mathbb{Z}_2 \times \mathbb{Z}_3$.

(b) $\{(m, n) \mid m=0 \text{ or } n=0, (m, n) \neq (0, 0)\}$

(c) Another example is $\mathbb{Z} \times \mathbb{Z}$. No, since by definition an integral domain D must contain the additive identity so we always have $(m, 0) \cdot (0, n) = (0, 0)$ in $D \times D$.

$$13. (a) \quad (a + b)(c + d) = (a + b)c + (a + b)d \\ = ac + bc + ad + bd$$

$$(b) \quad (a + b)(a + b) = aa + ba + ab + bb \quad \text{by part (a)} \\ = aa + ab + ab + bb \quad \text{since } R \text{ is commutative} \\ = a^2 + 2ab + b^2$$

15. Hint: The set of units of a ring is a group under multiplication. Apply a theorem from a group theory.

17. Proof of Corollary to Theorem 6.1.4: Since p is a prime, all nonzero elements of \mathbb{Z}_p are relatively prime to p . By Theorem 16.1.4 we are done.

Section 16.2

3. No, since $2^{-1} = 2$ in \mathbb{Z}_3 , but $a^{-1} \neq a$ and $b^{-1} \neq b$ in F .

5. (a) 0 (over \mathbb{Z}_2), 1 (over \mathbb{Z}_3), 3 (over \mathbb{Z}_5)

(b) 2 (over \mathbb{Z}_3), 3 (over \mathbb{Z}_5)

(c) 2

7. (a) 0 and 1 (b) 1 (c) 1 (d) none

9. (c) The roots of $x^2 - 2 = 0$ are $\sqrt{2}$ and $-\sqrt{2}$. Both numbers can be expressed in the form $a + b\sqrt{2}$ where $a, b \in \mathbb{Q}$: $\sqrt{2} = 0 + 1 \cdot \sqrt{2}$ and $-\sqrt{2} = 0 + -1 \cdot \sqrt{2}$.

(d) No, since $\pm\sqrt{3}$ cannot be expressed in the form $a + b\sqrt{2}$, $a, b \in \mathbb{Q}$. If there exist rational numbers a and b such that $\sqrt{3} = a + b\sqrt{2}$, then clearly $b \neq 0$ since $\sqrt{3}$ is irrational and $a \neq 0$ for that would imply that $\sqrt{3/2}$ is rational, which is false. If we square both sides, of the equation we will get a rational expression for $\sqrt{2}$ which is also false.

Section 16.3

1. (i) $f(x) + g(x) = 2 + 2x + x^2$, $f(x)g(x) = 1 + 2x + 2x^2 + x^3$

(ii) $f(x) + g(x) = x^2$, $f(x)g(x) = 1 + x^3$

(iii) $1 + 3x + 4x^2 + 3x^3 + x^4$

(iv) $1 + x + x^3 + x^4$

(v) $x^2 + x^3$

3. (a) If $a, b \in \mathbb{R}$, $a - b$ and ab are in \mathbb{R} since \mathbb{R} is a ring in its own right. Therefore, \mathbb{R} is a subring of $\mathbb{R}[x]$. The proofs of parts b and c are similar.

5. (a) Reducible, $(x + 1)(x^2 + x + 1)$

(b) Reducible, $x(x^2 + x + 1)$

(c) Irreducible. If you could factor this polynomial, one factor would be either x or $x + 1$, which would give you a root of 0 or 1, respectively. By substitution of 0 and 1 into this polynomial, it clearly has no roots.

(d) Reducible, $(x + 1)^4$

7. We illustrate this property of polynomials by showing that it is not true for a nonprime polynomial in $\mathbb{Z}_2[x]$. Suppose that $p(x) = x^2 + 1$, which can be reduced to $(x + 1)^2$, $a(x) = x^2 + x$, and $b(x) = x^3 + x^2$. Since $a(x)b(x) = x^5 + x^3 = x^3(x^2 + 1)$, $p(x) \mid a(x)b(x)$. However, $p(x)$ is not a factor of either $a(x)$ or $b(x)$.

9. The only possible proper factors of $x^2 - 3$ are $(x - \sqrt{3})$ and $(x + \sqrt{3})$, which are not in $\mathbb{Q}[x]$ but are in $\mathbb{R}[x]$.

11. For $n \geq 0$, let $S(n)$ be the proposition: For all $g(x) \neq 0$ and $f(x)$ with $\deg f(x) = n$, there exist unique polynomials $q(x)$ and $r(x)$ such that $f(x) = g(x)q(x) + r(x)$, and either $r(x) = 0$ or $\deg r(x) < \deg g(x)$.

Basis: $S(0)$ is true, for if $f(x)$ has degree 0, it is a nonzero constant, $f(x) = c \neq 0$, and so either $f(x) = g(x) \cdot 0 + c$ if $g(x)$ is not a constant, or $f(x) = g(x)g(x)^{-1} + 0$ if $g(x)$ is also a constant.

Induction: Assume that for some $n \geq 0$, $S(k)$ is true for all $k \leq n$. If $f(x)$ has degree $n + 1$, then there are two cases to consider. If $\deg g(x) > n + 1$, $f(x) = g(x) \cdot 0 + f(x)$, and we are done. Otherwise, if $\deg g(x) = m \leq n + 1$, we perform long division as follows, where LDT's = various terms of lower degree than $n + 1$.

$$g_m x^m + \text{LDT}' s \quad \frac{f_{n+1} \cdot g_m^{-1} x^{n+1-m}}{f_{n+1} x^{n+1} + \text{LDT}' s} \\ \frac{f_{n+1} x^{n+1} + \text{LDT}' s}{h(x)}$$

Therefore,

$$h(x) = f(x) - (f_{n+1} \cdot g_m^{-1} x^{n+1-m}) g(x) \Rightarrow \\ f(x) = (f_{n+1} \cdot g_m^{-1} x^{n+1-m}) g(x) + h(x)$$

Since $\deg h(x)$ is less than $n+1$, we can apply the induction hypothesis:

$$h(x) = g(x) q(x) + r(x) \text{ with } \deg r(x) < \deg g(x).$$

Therefore,

$$f(x) = g(x) (f_{n+1} \cdot g_m^{-1} x^{n+1-m} + q(x)) + r(x) \text{ with } \deg r(x) < \deg g(x).$$

This establishes the existence of a quotient and remainder. The uniqueness of $q(x)$ and $r(x)$ as stated in the theorem is proven as follows: if $f(x)$ is also equal to $g(x) \bar{q}(x) + \bar{r}(x)$ with $\deg \bar{r}(x) < \deg g(x)$, then

$$g(x) q(x) + r(x) = g(x) \bar{q}(x) + \bar{r}(x) \Rightarrow g(x) (\bar{q}(x) - q(x)) = r(x) - \bar{r}(x)$$

Since $\deg r(x) - \bar{r}(x) < \deg g(x)$, the degree of both sides of the last equation is less than $\deg g(x)$. Therefore, it must be that $\bar{q}(x) - q(x) = 0$, or $q(x) = \bar{q}(x)$ and so $r(x) = \bar{r}(x)$. ■

Section 16.4

1. If $a_0 + a_1 \sqrt{2} \in \mathbb{Q}[\sqrt{2}]$ is nonzero, then it has a multiplicative inverse:

$$\frac{1}{a_0 + a_1 \sqrt{2}} = \frac{1}{a_0 + a_1 \sqrt{2}} \frac{a_0 - a_1 \sqrt{2}}{a_0 - a_1 \sqrt{2}} \\ = \frac{a_0 - a_1 \sqrt{2}}{a_0^2 - 2a_1^2} \\ = \frac{a_0}{a_0^2 - 2a_1^2} - \frac{a_1}{a_0^2 - 2a_1^2} \sqrt{2}$$

The denominator, $a_0^2 - 2a_1^2$, is nonzero since $\sqrt{2}$ is irrational. Since $\frac{a_0}{a_0^2 - 2a_1^2}$ and $\frac{-a_1}{a_0^2 - 2a_1^2}$ are both rational numbers, $a_0 + a_1 \sqrt{2}$ is a unit of $\mathbb{Q}[\sqrt{2}]$. The field containing $\mathbb{Q}[\sqrt{2}]$ is denoted $\mathbb{Q}(\sqrt{2})$ and so $\mathbb{Q}(\sqrt{2}) = \mathbb{Q}[\sqrt{2}]$.

3. $x^4 - 5x^2 + 6 = (x^2 - 2)(x^2 - 3)$ has zeros $\pm\sqrt{2}$ and $\pm\sqrt{3}$. $\mathbb{Q}(\sqrt{2}) = \{a + b\sqrt{2} \mid a, b \in \mathbb{Q}\}$ contains the zeros $\pm\sqrt{2}$ but does not contain $\pm\sqrt{3}$, since neither are expressible in the form $a + b\sqrt{2}$. If we consider the set $\{c + d\sqrt{3} : c, d \in \mathbb{Q}(\sqrt{2})\}$, then this field contains $\pm\sqrt{3}$ as well as $\pm\sqrt{2}$, and is denoted $(\mathbb{Q}(\sqrt{2}))(\sqrt{3}) = \mathbb{Q}(\sqrt{2}, \sqrt{3})$. Taking into account the form of c and d in the description above, we can expand to

$$\mathbb{Q}(\sqrt{2}, \sqrt{3}) = \{b_0 + b_1\sqrt{2} + b_2\sqrt{3} + b_3\sqrt{6} \mid b_i \in \mathbb{Q}\}.$$

5. (a) $f(x) = x^3 + x + 1$ is reducible if and only if it has a factor of the form $x - a$. By Theorem 16.3.3, $x - a$ is a factor if and only if a is a zero. Neither 0 nor 1 is a zero of $f(x)$ over \mathbb{Z}_2 .

(b) Since $f(x)$ is irreducible over \mathbb{Z}_2 , all zeros of $f(x)$ must lie in an extension field of \mathbb{Z}_2 . Let c be a zero of $f(x)$. $\mathbb{Z}_2(c)$ can be described several different ways. One way is to note that since $c \in \mathbb{Z}_2(c)$, $c^n \in \mathbb{Z}_2(c)$ for all n . Therefore, $\mathbb{Z}_2(c)$ includes $0, c, c^2, c^3, \dots$. But $c^3 = c + 1$ since $f(c) = 0$. Furthermore, $c^4 = c^2 + c$, $c^5 = c^2 + c + 1$, $c^6 = c^2 + 1$, and $c^7 = 1$. Higher powers of c repeat preceding powers. Therefore,

$$\mathbb{Z}_2(c) = \{0, 1, c, c^2, c + 1, c^2 + 1, c^2 + c + 1, c^2 + c\} \\ = \{a_0 + a_1 c + a_2 c^2 \mid a_i \in \mathbb{Z}_2\}$$

The three zeros of $f(x)$ are c , c^2 and $c^2 + c$.

$$f(x) = (x + c)(x + c^2)(x + c^2 + c).$$

(c) Cite Theorem 16.2.4, part 3.

Section 16.5

3. Theorem 16.5.2 proves that not all nonzero elements in $F[[x]]$ are units.

$$\begin{aligned}
7. (a) \quad & b_0 = 1 \\
& b_1 = (-1)(2 \cdot 1) = -2 \\
& b_2 = (-1)(2 \cdot (-2) + 4 \cdot 1) = 0 \\
& b_3 = (-1)(2 \cdot 0 + 4 \cdot (-2) + 8 \cdot 1) = 0 \\
& \dots \text{ (all others are zero)}
\end{aligned}$$

$$\text{Hence, } f(x)^{-1} = 1 - 2x$$

$$\begin{aligned}
(b) \quad & f(x) = 1 + 2x + 2^2 x^2 + 2^3 x^3 + \dots \\
& = (2x)^0 + (2x)^1 + (2x)^2 + (2x)^3 + \dots \\
& = \frac{1}{1-2x}
\end{aligned}$$

The last step follows from the formula for the sum of a geometric series.

$$\begin{aligned}
9. (a) \quad & (x^4 - 2x^3 + x^2)^{-1} = (x^2(x^2 - 2x + 1))^{-1} \\
& = x^{-2}(1 - 2x + x^2)^{-1} \\
& = x^{-2} \left(\sum_{k=0}^{\infty} (k+1)x^k \right) \text{ by Example 2 of 16.5} \\
& = \sum_{k=-2}^{\infty} (k+2)x^k
\end{aligned}$$

Supplementary Exercises—Chapter 16

1. (a) This ring is not commutative.

$$\begin{aligned}
(A+B)^2 &= (A+B) \cdot (A+B) \\
&= (A+B) \cdot A + (A+B) \cdot B \\
&= A \cdot A + B \cdot A + A \cdot B + B \cdot B \\
&= A^2 + B \cdot A + A \cdot B + B^2
\end{aligned}$$

(b) Yes

3. (a) By Theorem 16.1.1 show:

(1) $[D, +]$ is a subgroup of the group $[M_{2 \times 2}(\mathbb{R}); +]$. We leave this to the reader.

(2) D is closed under multiplication. To prove this, let $\begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix}, \begin{pmatrix} c & 0 \\ 0 & d \end{pmatrix} \in D$. Then,

$$\begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix} \begin{pmatrix} c & 0 \\ 0 & d \end{pmatrix} = \begin{pmatrix} ac & 0 \\ 0 & bd \end{pmatrix} \in D$$

since ac and bd are real numbers and the product is in the form of a typical matrix in D .

(b) Since

$$\begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix} \begin{pmatrix} c & 0 \\ 0 & d \end{pmatrix} = \begin{pmatrix} ac & 0 \\ 0 & bd \end{pmatrix} = \begin{pmatrix} c & 0 \\ 0 & d \end{pmatrix} \begin{pmatrix} a & 0 \\ 0 & b \end{pmatrix},$$

D is commutative. The unity for D is $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$.

(c) The product of two nonzero matrices can be equal to zero. For example, $\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$. Therefore, D has divisors of zero and by Theorem 16.1.2 the cancellation law is not true in D .

5. (a) $2^4 = 16$

(b) The product cited in the solution to 3(c) above shows that $M_{2 \times 2}(\mathbb{R})$ has divisors of zero. Therefore, the matrix polynomial $(x - I)(x + I)$ may have solutions other than $\pm I$. In fact you can verify that $\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$ and $\begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}$ satisfy the given equation.

7. Use $T: A \rightarrow \mathbb{R}$ defined by $T\left(\begin{pmatrix} a & 0 \\ 0 & 0 \end{pmatrix}\right) = a$

9. By substitution and the operation tables of Example 16.2.2,

$$\begin{aligned}
a^2 + a + 1 &= b + a + 1 \\
&= 1 + 1 = 0
\end{aligned}$$

Therefore, a is a root. A similar calculation shows that b is a root. Substitution of 0 and 1 for x shows that they are not root.

11. By Theorem 16.3.3, $a \in \mathbb{Q}$ is a zero of $f(x)$ iff $(x - a)$ is a factor of $f(x)$, which also implies a must be a factor of 9. Hence, the only possible rational roots are: ± 1 , ± 3 , and ± 9 . We can verify that $(x - 3)$ is a divisor of $f(x)$ or that $x = 3$ is a zero of $f(x)$. Dividing $f(x)$ by $(x - 3)$ produces $q(x) = x^3 - 3x^2 + x - 3$, which has $x = 3$ as a rational root. Dividing $q(x)$ by $x - 3$ produces $x^2 + 1$. Hence, the complete factorization of $f(x)$ in $\mathbb{Q}[x]$ is $(x - 3)^2(x^2 + 1)$.

13. $g(0) = 0, g(1) = 1,$

$$g(a) = a^3 + a^2 + a = 1 + b + a = 1 + 1 = 0, \text{ and}$$

$$g(b) = b^3 + b^2 + b = 1 + a + b = 1 + 1 = 0.$$

Hence, 0, a , and b are zeros of $g(x)$ and the $g(x) = x(x - a)(x - b) = x(x + a)(x + b)$.

15. (a) Sum = (1, 0, 1), Product = (0, 1, 1, 1)

(b) Sum = (1, 0, 0, 0), Product = (0, 1, 1, 1, 0, 0, 1)

(c) Sum = (1, 1, 1, 0, 0), Product = (0, 0, 0, 0, 1, 1, 1, 0, 1)

(d) Sum = 010, Product = 11011

16. The encoding of a string of bits is based on polynomial division. Given a four bit message, we make the bits coefficients of a sixth degree polynomial, $b_3x^3 + b_4x^4 + b_5x^5 + b_6x^6$ which we can also express in \mathbb{Z}_2^6 as $(0, 0, 0, b_3, b_4, b_5, b_6)$, we divide this polynomial by $p(x) = 1 + x + x^3$ and add the remainder to the "message polynomial". The quotient in the division is discarded. Thus, if the remainder, which must be a polynomial of degree less than 2, is $b_0 + b_1x + b_2x^2$, the encoded message is the string of bits $(b_0, b_1, b_2, b_3, b_4, b_5, b_6)$.

(a) Encode the following elements of \mathbb{Z}_2^6 as described above.

(a) (0, 0, 0, 1, 1, 0, 1)

(b) (0, 0, 0, 1, 1, 1, 1)

(c) (0, 0, 0, 0, 0, 1, 0)

(b) Prove that the encoded message will always represent a polynomial which is evenly divisible by the polynomial $p(x)$ that is used to encode the message.

17. If the message polynomial is $m(x) = b_3x^3 + b_4x^4 + b_5x^5 + b_6x^6$ we divide by $p(x) = 1 + x + x^3$ and get a quotient and remainder: $m(x) = p(x)q(x) + r(x)$, where the degree of $r(x)$ is less than 3. We transmit $t(x) = m(x) + r(x) = m(x) + (m(x) - p(x)q(x)) = p(x)q(x)$ since $m(x) + m(x) = 0$. Now assume that the error x^k is added and we receive $p(x)q(x) + x^k$. Since $x^k, 0 \leq k \leq 6$, is not a multiple of $p(x)$, the received polynomial is also not a multiple of $p(x)$. The following *Mathematica* calculation verifies this last claim.

```
{x#, PolynomialRemainder[x#, x^3 + x + 1, x, Modulus -> 2]} & /@ Range[0, 6] //
Prepend[#, {"Monomial", "Remainder"}] &
```

Monomial	Remainder
1	1
x	x
x^2	x^2
x^3	$x + 1$
x^4	$x^2 + x$
x^5	$x^2 + x + 1$
x^6	$x^2 + 1$

19. (a) $b(x) = x^5 + x^4 + 1 = g(x)(x^2 + x + 1) + 0 \Rightarrow a = 111$

(b) $b(x) = x^5 + x^3 + x^2 + 1 = g(x)x^2 + 1$
 \Rightarrow error in the first bit of b
 $\Rightarrow e(a) = 001101$
 $\Rightarrow a = 001$

Getting a from $e(a)$ involves doing this calculation:

```
PolynomialQuotient[x^5 + x^3 + x^2, x^3 + x + 1, x, Modulus -> 2]
```

$$x^2$$

$$\begin{aligned}
 \text{(c) } b(x) &= x^5 + x + 1 = g(x)(x^2 + 1) + x^2 \\
 &\Rightarrow \text{error in the third bit of } b \\
 &\Rightarrow e(a) = 111001 \\
 &\Rightarrow a = 101
 \end{aligned}$$

PolynomialQuotient[$x^5 + x^2 + x + 1, x^3 + x + 1, x, \text{Modulus} \rightarrow 2$]

$$x^2 + 1$$

$$\begin{aligned}
 \text{(d) } b(x) &= x^4 + x^3 + x + 1 = g(x)(x + 1) + x^2 + x \\
 &\Rightarrow \text{error in the fifth bit of } b \\
 &\Rightarrow e(a) = 110100 \text{ (the string representation of } g(x)) \\
 &\Rightarrow a = 100
 \end{aligned}$$

21. (a) $g(x)$ is irreducible over \mathbb{Z}_2 since $g(0) = g(1) = 1$. Hence, $g(x)$ does not split in \mathbb{Z}_2 . Let β be a zero of $g(x)$, so that $\mathbb{Z}_2[\beta] = \{a + b\beta + c\beta^2 \mid a, b, c \in \mathbb{Z}_2\}$. This is a field of $2^3 = 8$ elements which, by Theorem 16.2.4, is isomorphic to $\text{GF}(8)$.

23. $1/g(x) = f(x)$ of Example 16.5.2.

$$\begin{array}{r}
 1 + 2x + 3x^2 + 4x^3 + \dots \\
 1 - 2x + x^2 \overline{)1} \\
 \underline{1 - 2x + 2x^2} \\
 2x - x^2 \\
 \underline{2x - 4x^2 + 2x^3} \\
 3x^2 - 2x^3 \\
 \underline{3x^2 - 6x^3 + 3x^4} \\
 4x^3 - 3x^4 \\
 \underline{4x^3 - 8x^4 + 4x^5} \\
 5x^4 - 4x^5 \\
 \vdots
 \end{array}$$

25. (a) $a_0 = a_1 = 1, a_2 = 2, a_3 = 3, a_4 = 5, \dots$, so

$$f(x) = 1 + x + 2x^2 + 3x^3 + 5x^4 + \dots$$

(b) $a_0 = a_1 = 1, a_2 = 0, a_3 = 1, a_4 = 1, a_5 = 0, \dots$

$$\begin{aligned}
 g(x) &= 1 + x + 0x^2 + x^3 + x^4 + 0x^5 + x^6 + x^7 + \dots \\
 &= (1 + x) + x^3(1 + x) + x^6(1 + x) + \dots \\
 &= (1 + x)(1 + x^3 + x^6 + \dots) \\
 &= \frac{(1+x)}{(1-x^3)}
 \end{aligned}$$

